

Reinforcement Learning for Decision-Making and Control in Power Systems: Tutorial, Review, and Vision

Xin Chen, Guannan Qu, Yujie Tang, Steven Low, Na Li

Abstract—With large-scale integration of renewable generation and distributed energy resources (DERs), modern power systems are confronted with new operational challenges, such as growing complexity, increasing uncertainty, and aggravating volatility. Meanwhile, more and more data are becoming available owing to the widespread deployment of smart meters, smart sensors, and upgraded communication networks. As a result, data-driven control techniques, especially reinforcement learning (RL), have attracted surging attention in recent years. In this paper, we provide a tutorial on various RL techniques and how they can be applied to decision-making and control in power systems. We illustrate RL-based models and solutions in three key applications, including frequency regulation, voltage control, and energy management. We conclude with three critical issues in the application of RL, i.e., safety, scalability, and data. Several potential future directions are discussed as well.

Index Terms—Reinforcement learning, smart grid, frequency regulation, voltage control, energy management.

NOMENCLATURE

A. Notations

\mathcal{S}, s	State space, state.
\mathcal{A}, a	Action space, action.
$\Delta(\mathcal{A})$	The set of probability distributions over set \mathcal{A} .
\mathbb{P}	Transition probability.
r	Reward.
o	Observation.
π, π^*	Policy, optimal policy.
γ	Discounting factor.
J	Expected total discounted reward.
Q_π	Q -function (or Q -value) under policy π .
$\text{NN}(x; w)$	Neural network with input x and parameter w .
\mathcal{N}	$:= \{1, \dots, N\}$, the set of buses in a power network or the set of agents.
\mathcal{E}	$\subseteq \mathcal{N} \times \mathcal{N}$, the set of lines connecting buses.
\mathcal{T}	$:= \{0, 1, \dots, T\}$, the discrete time horizon.
Δt	The time interval in \mathcal{T} .

B. Abbreviations

A3C	Asynchronous Advantage Actor Critic.
AMI	Advanced Metering Infrastructure.

X. Chen, Y. Tang, and N. Li are with the School of Engineering and Applied Sciences, Harvard University, USA. Email: chen_xin@g.harvard.edu, yujietang@seas.harvard.edu, nali@seas.harvard.edu.

G. Qu and S. Low are with Department of Computing and Mathematical Sciences, California Institute of Technology, USA. Email: gqu@caltech.edu, slow@caltech.edu.

ANN	Artificial Neural Network.
DDPG	Deep Deterministic Policy Gradient.
DER	Distributed Energy Resource.
DP	Dynamic Programming.
(D)RL	(Deep) Reinforcement Learning.
DQN	Deep Q Network.
EMS	Energy Management System.
EV	Electric Vehicle.
FR	Frequency Regulation.
HVAC	Heating, Ventilation, and Air Conditioning.
IES	Integrated Energy System.
LSPI	Least-Squares Policy Iteration.
LSTM	Long-Short Term Memory.
MDP	Markov Decision Process.
OLTC	On-Load Tap Changing Transformer.
OPF	Optimal Power Flow.
PMU	Phasor Measurement Unit.
SAC	Soft Actor Critic.
SCADA	Supervisory Control and Data Acquisition.
SVC	Static Var (Reactive Power) Compensator.
TD	Temporal Difference.
UCRL	Upper Confidence Reinforcement Learning.

I. INTRODUCTION

ELECTRIC power systems are undergoing an architectural transformation to become more sustainable, distributed, dynamic, intelligent, and open. On the one hand, the proliferation of renewable generation and distributed energy resources (DERs), including solar energy, wind power, energy storage, responsive demands, electric vehicles (EVs), etc., creates severe operational challenges. On the other hand, the deployment of information, communication, and computing technologies throughout the electric system, such as phasor measurement units (PMUs), advanced metering infrastructures (AMIs), and wide area monitoring systems (WAMS) [1], has been growing rapidly in recent decades. It evolves traditional power systems towards smart grids and offers an unprecedented opportunity to overcome these challenges through real-time data-driven monitoring and control at scale. This will require new advanced decision-making and control techniques to manage:

- 1) *Growing complexity*. The deployment of massive DERs and interconnection of regional power grids dramatically increase the complexity of system operation and make it difficult to obtain accurate system (dynamical) models.
- 2) *Increasing uncertainty*. The rapid growth of renewable generation and responsive loads greatly increases uncer-

tainty, especially when human users are involved, which jeopardizes predictions and system reliability.

- 3) *Aggravating volatility*. The high penetration of power electronics converter-interfaced devices reduces system inertia, which leads to faster dynamics and necessitates advanced controllers with online adaptivity.

In particular, reinforcement learning (RL) [2], a prominent machine learning paradigm that is concerned with how agents take sequential actions in an uncertain interactive environment and learn from the feedback to optimize a certain performance, can play an important role in overcoming these challenges. Leveraging artificial neural networks (ANNs) for function approximation, deep RL (DRL) [3] is further developed to solve large-scale online decision problems. The most appealing virtue of (D)RL is its *model-free* nature, i.e., it makes decisions without explicitly estimating the underlying models. Therefore (D)RL has the potential to capture hard-to-model dynamics and could outperform model-based methods in highly complex tasks. Moreover, the data-driven nature of (D)RL allows it to adapt to real-time observations and perform well in uncertain dynamical environments. The past decade has witnessed great success of (D)RL in a broad spectrum of applications, such as playing games [4], robotics [5], autonomous driving [6], clinical trials [7], and etc.

Meanwhile, the application of RL in power system operation and control has attracted surging attention [8]–[11]. The RL-based decision-making mechanisms are envisioned to compensate for the limitations of existing model-based approaches, and thus are promising to address the emerging challenges described above. In this paper, we provide a tutorial and review on RL-based decision-making in power systems. We will introduce various RL terminologies, elaborate on how RL can be applied to power systems, and discuss critical issues in their application. Comparing with recent review articles [8]–[11] on this subject, the main merits of this paper include

- 1) A comprehensive overview of RL methodology is presented, from basic concepts and theoretical fundamentals to state-of-the-art RL techniques.
- 2) Rather than listing the relevant literature for broad power system applications, we select three key applications to illustrate how to model and solve with RL methods in detail and with mathematical models.
- 3) The key limitations and future directions of applying RL to power systems are discussed in depth.

In the remainder of this paper, Section II presents a comprehensive overview of the RL fundamentals and its state-of-the-art techniques. Section III describes the application of RL to three critical power system problems, i.e., *frequency regulation*, *voltage control*, and *energy management*. Paradigmatic mathematical models are provided for illustration as well. Section IV summarizes three key issues in these applications, i.e., safety, scalability and data, and discusses several potential future directions. Lastly, we draw conclusions in Section V.

II. PRELIMINARIES ON REINFORCEMENT LEARNING

This section provides a comprehensive overview on the RL methodology. Firstly, we set up the RL problem formulation

and some key concepts, such as Q -function and Bellman (Optimality) Equation. Then two categories of classical RL algorithms, i.e., value-based and policy-based, are introduced. With these fundamentals in place, we next present several state-of-the-art RL techniques, involving DRL, deterministic policy gradient, modern actor-critic methods, multi-agent RL, and etc. The overall structure of RL methodology with related literature is illustrated as Figure 1.

A. Fundamentals of Reinforcement Learning

RL is a branch of machine learning concerned with how an agent makes sequential decisions in an uncertain **environment** to maximize the cumulative reward. Mathematically, the decision-making problem is modeled as Markov Decision Processes (MDPs), which are defined by **state** space \mathcal{S} , **action** space \mathcal{A} , the **transition** probability function $\mathbb{P}(\cdot|s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ that maps a state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ to a distribution on the state space, and lastly the **reward** function $r(s, a)^1 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The state space \mathcal{S} and action space \mathcal{A} can be either discrete or continuous. To simplify discussion, we focus on the discrete case in the follows.

As illustrated in Figure 2, in a MDP problem, the environment starts with an initial state $s_0 \in \mathcal{S}$. At each time $t = \{0, 1, \dots\}$, given current state $s_t \in \mathcal{S}$, the agent chooses action $a_t \in \mathcal{A}$ and receives reward $r(s_t, a_t)$ that depends on the current state-action pair (s_t, a_t) , after which the next state s_{t+1} is randomly generated from the transition probability $\mathbb{P}(s_{t+1}|s_t, a_t)$. A **policy** $\pi(a|s) \in \Delta(\mathcal{A})$ for the agent is a map from the state s to a distribution on the action space \mathcal{A} , which provides a rule on what action to take given a certain state s .² The agent aims to find an optimal policy π^* (may not be unique) that maximizes the expected infinite horizon discounted reward $J(\pi)$:

$$\pi^* \in \arg \max_{\pi} J(\pi) = \mathbb{E}_{s_0 \sim \mu_0} \mathbb{E}_{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t), \quad (1)$$

where the first expectation means that s_0 is drawn from an initial state distribution μ_0 , and the second expectation means that the action a_t is taken according to policy $\pi(\cdot|s_t)$. Parameter $\gamma \in (0, 1)$ is the discounting factor that penalizes the rewards in the future.

In the MDP framework, the so-called “**model**” specifically refers to the reward function r and the transition probability \mathbb{P} . Accordingly, it leads to two different problem settings:

- *When the model is known*, one can directly solve for an optimal policy π^* by Dynamic Programming (DP) [13].
- *When the model is unknown*, the agent learns an optimal policy π^* based on the past observations from interacting with the environment, which is the problem of RL.

Since DP lays the foundation for RL algorithms, we first consider the case with known model and introduce the basic

¹A generic reward function is given by $r(s, a, s')$ where the next state s' is also included as an argument, but there is no essential difference between the case with $r(s, a)$ and the case with $r(s, a, s')$ in algorithms and results. By marginalizing over next states s' according to the transition function $\mathbb{P}(s'|s, a)$, one can simply convert $r(s, a, s')$ to $r(s, a)$ [2].

² $a \sim \pi(\cdot|s)$ is a stochastic policy, and it becomes a deterministic policy $a = \pi(s)$ when the probability distribution $\pi(\cdot|s)$ is a singleton for all s .

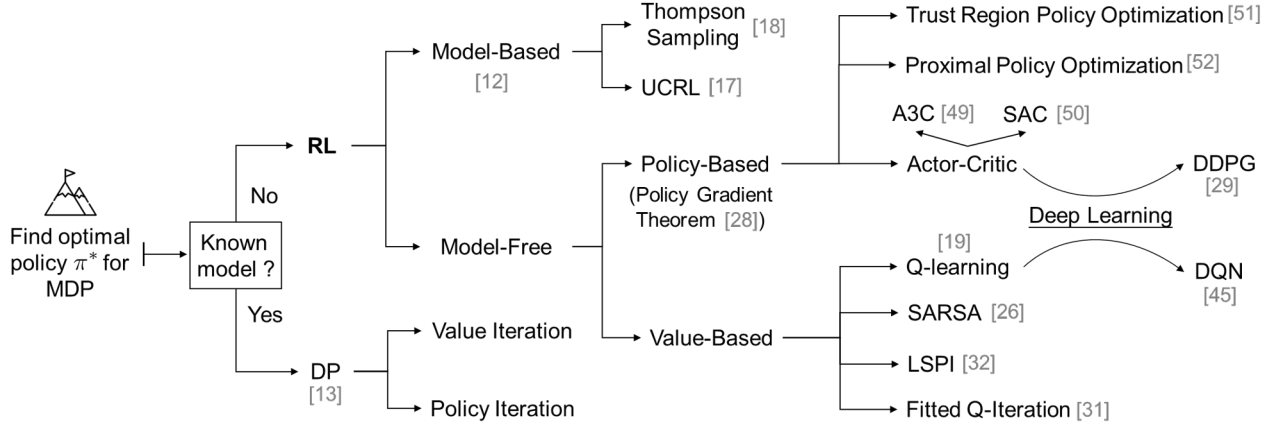


Fig. 1. Structure of the RL methodology with related literature. (“Model-based” refers to the RL algorithms that explicitly estimate and online update a system model and take actions based on this estimated model [12]. In contrast, “model-free” means that the associated RL algorithms directly search for optimal policies based on Q -function or policy gradient methods, without estimating the system model.)

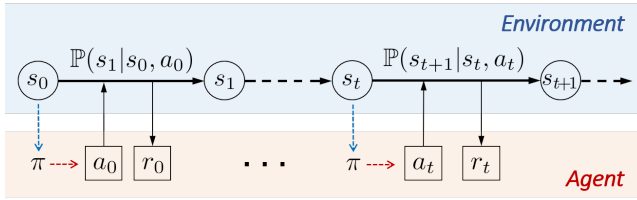


Fig. 2. Illustration of a Markov Decision Process.

ideas of finding an optimal policy π^* with respect to (1). The crux is the concept of Q -function together with the *Bellman Equation*. The Q -function $Q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ for a given policy π is defined as

$$Q_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right], \quad (2)$$

which is the expected cumulative reward when the initial state is s , the initial action is a , and all the subsequent actions are chosen according to policy π . The Q -function Q_π satisfies the following *Bellman Equation*: $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$,

$$Q_\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a), a' \sim \pi(\cdot | s')} Q_\pi(s', a'), \quad (3)$$

where the expectation denotes that the next state s' is drawn from $\mathbb{P}(\cdot | s, a)$, and the next action a' is drawn from $\pi(\cdot | s')$. Here, it is helpful to think of the Q -function as a large table or vector filled with Q -values $Q_\pi(s, a)$. The *Bellman Equation* (3) indicates a recursive relation that each Q -value equals to the immediate reward plus the discounted future value. Computing the Q -function for a given policy π is called *policy evaluation*, which can be done by simply solving a set of linear equations when the model, i.e., \mathbb{P} and r , is known.

The Q -function corresponding to an optimal policy π^* for (1) is called an *optimal Q -function* and denoted as Q^* . The key to find π^* is that the optimal Q -function must be the unique solution to the *Bellman Optimality Equation* (4): $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$,

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \max_{a' \in \mathcal{A}} Q^*(s', a'). \quad (4)$$

Interested readers are referred to textbook [13, Sec. 1.2] on why this is true. Based on the *Bellman Optimality Equation* (4), the optimal Q -function Q^* and optimal policy π^* can be solved using DP or linear programming [13, Sec. 2.4]. Two classic DP algorithms are *policy iteration* and *value iteration*, and see [2, Chapter 4] for details.

Remark 1. (Modeling Issues with MDP). MDP is a generic framework to model sequential decision-making problems and is the basis for RL algorithms. However, when modeling power system control problems in the MDP framework, there are several issues that deserve attention.

- 1) At the heart of MDP is the *Markov property* that the distribution of future states depends only on the present state and action, i.e., $\mathbb{P}(s_{t+1} | s_t, a_t)$. In other words, given the present, the future does not depend on the past. Then for a specific control problem, it needs to check whether the choices of state and action satisfy the Markov property. A general guideline is to include all necessary known information in the enlarged state, known as *state augmentation* [13], to maintain the Markov property, however at the cost of complexity.
- 2) A majority of classical MDP theories and RL algorithms are based on *discrete-time transitions*. While many power system control problems follow continuous-time dynamics, such as frequency regulation. To fit the MDP framework, continuous-time dynamics are usually discretized with a proper temporal resolution, which is a common issue for most digital control systems and there are well-established frameworks to deal with it. Nevertheless, there are RL variants that are built directly on continuous-time dynamics, such as *integral RL* [14].
- 3) Many MDP/RL methods assume time-homogeneous state transitions and rewards. However, there are various time-varying exogenous inputs and disturbances in power systems, making the state transitions *not time-homogeneous*. This is an important problem that has not been adequately explored in the existing power literature and deserves further study. *Nonstationary MDP* [15] and related RL

algorithms [16] could be the potential directions.

Besides, the issues of continuous state/action spaces and partial observability for MDP modeling will be discussed later.

B. Classical Reinforcement Learning Algorithms

This subsection considers the RL setting when the environment model is unknown, and presents classical RL algorithms for finding the optimal policy π^* . Model-free RL algorithms³ are mainly categorized into two types: *value-based* and *policy-based*. Generally, for modest-scale RL problems with finite state/action space, value-based methods are preferred as they do not assume a policy class and have strong convergence guarantee. The convergence of value-based methods to an optimal Q -function in the tabular setting (without function approximation) was proven back in the 1990s [19]. In contrast, policy-based methods are more efficient for problems with high dimension or continuous action/state space. But they are known to suffer from various convergence issues, e.g., local optimum, high variance, etc., and little is known about their basic theoretical convergence properties. Until recently, convergence of policy-based methods with restricted policy classes to the global optimum under tabular policy parameterization has been proven [20].

Remark 2. (Exploration vs. Exploitation). A fundamental problem faced by both value-based and policy-based RL algorithms is the dilemma between exploration and exploitation. Good performance requires taking actions in an adaptive way that strikes an effective balance between 1) *exploring* poorly-understood actions to gather new information that may improve future reward and 2) *exploiting* what is known for decisions to maximize immediate reward. Generally, it is natural to achieve exploitation with the goal of reward maximization, while different RL algorithms encourage exploration in different ways. For value-based RL algorithms, ϵ -greedy is commonly utilized with a probability of ϵ to explore random actions. In policy-based methods, the exploration is usually realized by adding random perturbation to the actions or adopting a stochastic policy.

Before presenting the two types of RL algorithms, we introduce a key algorithm, **Temporal-Difference (TD)** learning [21], for policy evaluation without knowing the model. TD learning is central to both value-based and policy-based RL algorithms. It learns the Q -function Q_π for a given policy π from episodes of experience. Here, an “episode” refers to a state-action-reward trajectory over time $(s_0, a_0, r_0, s_1, a_1, r_1, \dots)$ until terminated. Specifically, TD learning maintains a Q -function $Q(s, a)$ for all state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$ and updates it upon a new observation (r_t, s_{t+1}, a_{t+1}) by

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)), \quad (5)$$

where α is the step size. Readers might observe that the second term in (5) is very similar to the Bellman Equation (3), which

is exactly the rationale behind TD learning. Essentially, TD learning (5) is a stochastic approximation scheme for solving the Bellman Equation (3) [22], and can be shown to converge to the true Q_π under mild assumptions [21], [23].

1) *Value-based RL algorithms* directly learn the optimal Q -function Q^* , whilst the optimal (deterministic) policy π^* is a byproduct that can be retrieved by acting greedily, i.e., $\pi^*(s) = \arg \max_{a \in \mathcal{A}} Q^*(s, a)$. Among many, **Q -learning** [19], [24] is perhaps the most popular value-based RL algorithm. Similar to TD-learning, Q -learning maintains a Q -function and updates it towards the optimal Q -function based on episodes of experience. Specifically, at each time t , given current state s_t , the agent chooses action a_t according to a certain behavior policy.⁴ Upon observing the outcome (r_t, s_{t+1}) , Q -learning updates the Q -function by

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a' \in \mathcal{A}} Q(s_{t+1}, a') - Q(s_t, a_t)). \quad (6)$$

The rationale behind is that Q -learning algorithm (6) is essentially a stochastic approximation scheme for solving the Bellman Optimality Equation (4), and it can be shown to converge to Q^* under mild assumptions [19], [25].

SARSA⁵ [26] is another classical value-based RL algorithm, whose name comes from the experience sequence (s, a, r, s', a') . SARSA is actually an *on-policy* variant of Q -learning. The major difference is that SARSA takes actions according to the target policy (typically ϵ -greedy based on the current Q -function) rather than any arbitrary behavior policy in Q -learning. The following remark distinguishes and compares “on-policy” and “off-policy” RL algorithms.

Remark 3. (On-Policy vs. Off-Policy). On-policy RL methods continuously improve a policy (called target policy) and implement this policy to generate episodes for algorithm training. In contrast, off-policy RL methods learn a target policy based on the episodes generated by a different policy (called behavior policy) rather than the target policy itself. In short, “on” and “off” indicate whether the training samples are generated by following the target policy or not. For example, Q -learning is an off-policy method as the episodes used in training can be produced by any policies, while the actor-critic algorithm described below is on-policy.⁶ For power system applications, control policies that are not well trained are generally not allowed to be implemented in real-world power grids for the sake of safety. Thus off-policy RL is usually preferred when high-fidelity simulators are unavailable, since it can learn from the huge amount of operational data generated by incumbent controllers. Off-policy RL is also relatively easy to provide safety guarantee due to the flexibility in choosing the behavior policies, while it is known to suffer from slower convergence and higher sample complexity.

2) *Policy-based RL algorithms* restrict the optimal policy search to a policy class that is parameterized as π_θ with

⁴Such a behavior policy, also called exploratory policy, can be arbitrary as long as it visits all the states and actions sufficiently often.

⁵In some literature, SARSA is also called Q -learning.

⁶There are also off-policy variants, e.g., [27], of the actor-critic algorithm.

³There are also model-based RL algorithms, e.g., Upper Confidence RL [17] and Thompson sampling [18], which online estimate the environment model from past observations and make decisions based on the model.

parameter $\theta \in \Theta \subset \mathbb{R}^K$. With this parameterization, the objective (1) can be rewritten as a function of the policy parameter, i.e., $J(\theta)$, and the RL problem is reformulated as an optimization problem (7) to find the optimal θ^* :

$$\theta^* \in \arg \max_{\theta \in \Theta} J(\theta). \quad (7)$$

To solve (7), a straightforward idea is to employ the gradient ascent method, i.e., $\theta \leftarrow \theta + \eta \nabla J(\theta)$, where η is the step size. However, computing the gradient $\nabla J(\theta)$ was supposed to be intrinsically hard as the environment model is unknown. *Policy Gradient Theorem* [28] is a big breakthrough in addressing the gradient computation issue. This theorem shows that the policy gradient $\nabla J(\theta)$ can be simply expressed as

$$\nabla J(\theta) = \sum_{s \in \mathcal{S}} \mu_\theta(s) \sum_{a \in \mathcal{A}} \pi_\theta(a|s) Q_{\pi_\theta}(s, a) \nabla_\theta \ln \pi_\theta(a|s). \quad (8)$$

Here, $\mu_\theta(s) \in \Delta(\mathcal{S})$ is the on-policy state distribution [2, Chapter 9.2], which denotes the fraction of time steps spent in each state $s \in \mathcal{S}$. Expression (8) is for a stochastic policy $a \sim \pi_\theta(\cdot|s)$, while the version of policy gradient theorem for a deterministic policy $a = \pi_\theta(s)$ [29] will be discussed later.

The policy gradient theorem provides a highway to estimate the gradient $\nabla J(\theta)$, which lays the foundation for policy-based RL algorithms. In particular, the **actor-critic** algorithm is a prominent and widely used architecture based on policy gradient. It consists of two eponymous components: 1) the “critic” is in charge of estimating the Q -function $Q_{\pi_\theta}(s, a)$, and 2) the “actor” conducts the gradient ascent step based on (8). An illustrative actor-critic example is given by the following iterative scheme:

- 1) Given state s , take action $a \sim \pi_\theta(a|s)$, then observe the reward r and next state s' ;
- 2) (Critic) Update Q -function $Q_{\pi_\theta}(s, a)$ by TD learning;
- 3) (Actor) Update policy parameter θ by

$$\theta \leftarrow \theta + \eta Q_{\pi_\theta}(s, a) \nabla_\theta \ln \pi_\theta(a|s); \quad (9)$$

- 4) $s \leftarrow s'$. Go to step 1) and repeat.

Note that there are many variants of the actor-critic method with different implementation details, e.g., how s is sampled from $\mu_\theta(s)$, how Q -function is updated, etc. See [2, Chapter 13] for detailed introduction.

We emphasize that the algorithms introduced above are far from complete. In the next subsections, we will introduce the state-of-the-art modern (D)RL techniques that are widely used in complex control tasks, especially for power system applications. At last, we close this subsection with the following two remarks on different RL settings.

Remark 4. (Online RL vs. Batch RL). The algorithms introduced above are referred as “*online RL*” that takes actions and updates the policies simultaneously. In contrast, there is another type of RL called “*batch RL*” [30], which decouples the sample data collection and policy training. Specifically, given a set of experience episodes generated by any arbitrary behavior policies, batch RL fits the optimal Q -function or optimizes the target policy fully based on this fixed sample dataset. Some classical batch RL algorithms include Fitted Q -Iteration [31], Least-Squares Policy Iteration (LSPI) [32],

etc. For example, given a batch of transition experiences $\mathcal{D} := \{(s_i, a_i, r_i, s'_i)_{i=1}^n\}$, Fitted Q -Iteration, which is seen as the batch version of Q -learning, aims to fit a parameterized Q -function $Q_\theta(s, a)$ by iterating the following two steps:

- 1) Create the target Q -value q_i for each sample in \mathcal{D} by

$$q_i = r_i + \gamma \max_{a'} Q_\theta(s'_i, a').$$

- 2) Apply regression approaches to fit a new $Q_\theta(s, a)$ based on the training dataset $(s_i, a_i; q_i)_{i=1}^n$.

The crucial advantages of batch RL lie in the stability and data-efficiency of the learning process by making the best use of the available sample datasets. However, because of relying entirely on a given dataset, the lack of exploration is one of the major problems of batch RL. To encourage exploration, batch RL typically iterates between exploratory sample collection and policy learning prior to application. Besides, pure batch RL, also referred as “*offline RL*” [33], has attracted increasing recent attention, which completely ignores the exploration issue and aims to learn policies fully based on a static dataset without any online interaction. Offline RL generally assumes a sufficiently large and diverse dataset that adequately covers high-reward transitions for learning good policies, and turns the RL problem into a supervised machine learning problem. See [34] for a tutorial of offline RL.

Remark 5. (Passive RL, Active RL, and Inverse RL). In literature, the terminology “*passive RL*” typically refers to the RL setting where the agent acts based on a fixed policy π and aims to learn how good this policy is from observations. It is analogous to the policy evaluation task, and TD learning is one of the representative algorithms of passive RL. In contrast, “*active RL*” allows the agent to update policies with the goal of finding an optimal policy, which is basically the standard RL setting that we described above. However, in some references, e.g., [35], [36], “*active RL*” has a completely different meaning and refers to the RL variant where the agent does not observe the reward unless it pays a query cost, to account for the difficulty of collecting reward feedback. Thus at each time, the agent chooses both an action and whether to observe the reward. Another interesting RL variant is “*inverse RL*” [37], [38], in which the state-action sequence of an (expert) agent is given and the task is to infer the reward function that this agent seeks to maximize. Inverse RL is motivated by many practical applications where the reward engineering is complex or expensive while one can observe an expert demonstrating the task to learn how to perform, e.g., autonomous driving.

C. Fundamentals of Deep Learning

This subsection presents the fundamentals of deep learning to set the stage for the introduction of DRL. Deep learning refers to the machine learning technique that models with multi-layer ANNs. The history of ANN dates back to 1940s [39], and it has received tremendous interests in the recent decade due to the booming of data technology and computing power, which allows efficient training of wider and deeper ANNs. Essentially, ANN is an universal parameterized mapping $y = \text{NN}(x; w)$ from the input features x to the

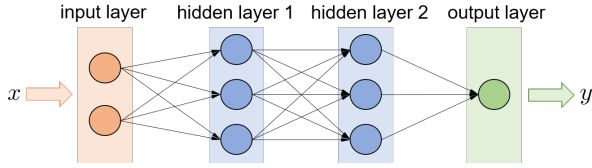


Fig. 3. Illustration of a regular four-layer feed-forward ANN.

outputs y with the parameters w . As illustrated in Figure 3, an input feature vector x is taken in by the input layer, then is processed through a series of hidden layers, and becomes the output vector y . Each hidden layer consists of a number of neurons that are the activation functions, e.g. linear, ReLU, or sigmoid [40]. Based on a sample dataset $(x^i, y^i)_{i=1, \dots, n}$, the parameter w can be optimized via regression. A landmark in training ANNs is the discovery of the *back-propagation* method, which offers an efficient way to compute the gradient of the loss function over w [41]. Nevertheless, it is pretty tricky to train large-scale ANNs in practice, and article [41] provides an overview of the optimization algorithms and theory for training ANNs. Three typical classes of ANNs with different architectures are introduced below. See book [42] for details.

1) *Convolutional Neural Networks (CNNs)* are in the architecture of feed-forward neural networks (as shown in Figure 3) and specialize in pattern detection, which are powerful for image analysis and computer vision tasks. The convolutional hidden layers are the basis at the heart of a CNN. Each neuron $k = 1, 2, \dots$ in a convolutional layer defines a small filter (or kernel) matrix F_k of low dimension (e.g., 3×3) and convolves with the input matrix X of relatively high dimension, which leads to the output matrix $U_k = F_k \otimes X$. Here, \otimes denotes the convolution operator⁷, and the output $(U_k)_{k=1,2,\dots}$ is referred as the feature map that is passed to the next layer. Besides, pooling layers are commonly used to reduce the dimension of the representation with the max or average pooling.

2) *Recurrent Neural Networks (RNNs)* specialize in processing long sequential inputs and tackling tasks with context spreading over time by leveraging a recurrent structure. Hence RNNs achieve great success in the applications such as speech recognition and machine translation. RNNs process an input sequence one element at a time, and maintain in their hidden units a state vector s that implicitly contains historical information about the past elements. Interestingly, the recurrence of RNNs is analogous to a dynamical system [42] and can be expressed as

$$s_t = f(s_{t-1}, x_t; v), \quad y_t = g(s_t; u), \quad (10)$$

where x_t and y_t are the input and output of the neural network at time step t , and $w := (v, u)$ is the parameter for training. s_t denotes the state stored in the hidden units at step t and will be passed to the processing at step $t+1$. In this way, s_t implicitly covers all historical input information (x_1, \dots, x_t) . Among

⁷Specifically, the convolution operation is performed by sliding the filter matrix F_k across the input matrix X and computing the corresponding element-wise multiplications, so that each element in matrix U_k is the sum of the element-wise multiplications between F_k and the associated sub-matrix of X . See [42, Chapter 9] for a detailed definition of convolution.

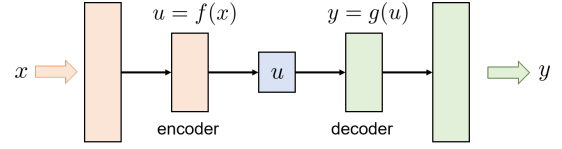


Fig. 4. Illustration of a simple autoencoder network.

many variants, long-short term memory (LSTM) network [43] is a special type of RNNs that excels at handling long-term dependencies and outperforms conventional RNNs by using special memory cells and gates.

3) *Autoencoders* [44] are used to obtain a low-dimensional representation of high-dimensional inputs, which is similar to, but more general than, principal components analysis (PCA). As illustrated in Figure 4, an autoencoder is in an hourglass-shape feed-forward network structure and consists of an encoder function $u = f(x)$ and a decoder function $y = g(u)$. In particular, an autoencoder is trained to learn an approximation function $y = \text{NN}(x; w) = g(f(x)) \approx x$ with the loss function $L(x, g(f(x)))$ that penalizes the dissimilarity between the input x and output y . The bottleneck layer has a much smaller amount of neurons, thus it is forced to form a compressed representation of the input x .

D. Deep Reinforcement Learning

For many practical problems, the state and action spaces are large or continuous, together with complex system dynamics. As a result, it is intractable for value-based RL to compute or store a gigantic Q -value table for all state-action pairs. To deal with this issue, *function approximation* methods are developed to approximate the Q -function with some parameterized function classes, such as linear function or polynomial function. As for policy-based RL, finding a capable policy class to achieve optimal control is also nontrivial in high-dimensional complex tasks. Driven by the advances of deep learning, DRL that leverages ANNs for function approximation or policy parameterization is becoming increasingly popular. Specifically, DRL can use ANNs to 1) approximate the Q -function with a Q -network $\hat{Q}_w(s, a) := \text{NN}(s, a; w)$, and 2) parameterize the policy with the policy network $\pi_\theta(a|s) := \text{NN}(a|s; \theta)$. We elaborate these two usages as follows.

1) *Q-Function Approximation*. Q -network can be used to approximate the Q -function in TD learning (5) and Q -learning (6). For TD learning, the parameter w is updated by

$$w \leftarrow w + \alpha \left[r_t + \gamma \hat{Q}_w(s_{t+1}, a_{t+1}) - \hat{Q}_w(s_t, a_t) \right] \nabla_w \hat{Q}_w(s_t, a_t), \quad (11)$$

where the gradient $\nabla_w \hat{Q}_w(s_t, a_t)$ can be calculated efficiently using the back-propagation method. As for Q -learning, however, it is known that adopting a nonlinear function, such as a ANN, for approximation may cause instability and divergence issues in the training process. To this end, **Deep Q-Network (DQN)** [45] is developed and greatly improves the training stability of Q -learning with the following two tricks:

- *Experience Replay*. Instead of performing on consecutive episodes, a commonly used trick is to store all the transition experiences $e := (s, a, r, s')$ in a database \mathcal{D} called “replay buffer”. At each step, a batch of transition experiences is randomly sampled from the replay buffer \mathcal{D} for Q -learning update. This can enhance the data efficiency by recycling previous experiences and reduce the variance of learning updates. More importantly, sampling uniformly from the replay buffer breaks the temporal correlations that jeopardize the training process, and thus improves the stability and convergence of Q -learning.

- *Target Network*. The other trick is the introduction of the target network $\hat{Q}_{\hat{w}}(s, a)$ with parameter \hat{w} , which is a clone of the Q -network $\hat{Q}_w(s, a)$, while its parameter \hat{w} is kept frozen and only gets updated periodically. Specifically, with a batch of transition experiences $(s_i, a_i, r_i, s'_i)_{i=1}^n$ sampled from the replay buffer, the Q -network $\hat{Q}_w(s, a)$ is updated by solving

$$w \leftarrow \arg \min_w \sum_{i=1}^n (r_i + \gamma \max_{a'} \hat{Q}_{\hat{w}}(s'_i, a') - \hat{Q}_w(s_i, a_i))^2. \quad (12)$$

The optimization (12) can be viewed as finding an optimal Q -network $\hat{Q}_w(s, a)$ that approximately solves the Bellman Optimality Equation (4). Whilst the critical difference is that the target network $\hat{Q}_{\hat{w}}$ with parameter \hat{w} instead of \hat{Q}_w is used to compute the maximization over a' in (12). After a fixed number of updates above, the target network $\hat{Q}_{\hat{w}}(s, a)$ is renewed by replacing \hat{w} with the latest learned w . This trick can mitigate the training instability as the short-term oscillations are circumvented. See [46] for more details.

In addition, there are several notable variants of DQN that further improve the performance, such as double DQN [47] and dueling DQN [48]. Specifically, double DQN is proposed to tackle the overestimation issue of action values in DQN by learning two sets of Q -functions; one Q -function is used to select the action and the other is used to determine its value. Dueling DQN proposes a dueling network architecture that separately estimates the state value function $V(s)$ and the state-dependent action advantage function $A(s, a)$, which are then combined to determine the Q -value. The main benefit of this factoring is to generalize learning across actions without imposing any change to the underlying RL algorithm [48].

2) *Policy Parameterization*. Due to the powerful generalization capability, ANNs are widely used to parameterize control policies, especially when the state and action spaces are continuous. The resultant policy network $\text{NN}(a|s; \theta)$ takes states as the input and outputs the probability of action selection. In actor-critic methods, it is common to adopt both the Q -network $\text{NN}(s, a; w)$ and the policy network $\text{NN}(a|s; \theta)$ simultaneously, where the “actor” updates θ according to (9) and the “critic” updates w according to (11). The back-propagation method [41] can be used to compute the gradient of ANNs efficiently.

When function approximation is adopted, the theoretical analysis on both value-based and policy-based RL methods is little and generally limited to the linear function approximation. Besides, one problem that hinders the use of value-based methods for large or continuous action space is the difficulty

of performing the maximization step. For example, when deep ANNs are used to approximate the Q -function, it is not easy to solve $\max_{a'} \hat{Q}_w(s, a')$ for the optimal action a' due to the nonlinear complex formulation of $\hat{Q}_w(s, a)$.

E. Other Modern Reinforcement Learning Techniques

This subsection summarizes several state-of-the-art modern RL techniques that are widely used in complex tasks.

1) *Deterministic Policy Gradient*: The RL algorithms described above focus on stochastic policies $a \sim \pi_\theta(\cdot|s)$, while deterministic policies $a = \pi_\theta(s)$ are more desirable for many real-world control problems with continuous state and action spaces. On the one hand, since most incumbent controllers in physical systems, such as PID control and robust control, are all deterministic, deterministic policies are better matched to the practical control architectures, e.g., in power system applications. On the other hand, a deterministic policy is more sample-efficient as its policy gradient only integrates over the state space, whereas a stochastic policy gradient integrates over both state and action spaces [29]. Similar to the stochastic case, there is *Deterministic Policy Gradient Theorem* [29] showing that the policy gradient for a deterministic policy $\pi_\theta(s)$ can be simply expressed as

$$\nabla J(\theta) = \mathbb{E}_{s \sim \mu_\theta(\cdot)} \nabla_a Q_{\pi_\theta}(s, a)|_{a=\pi_\theta(s)} \nabla_\theta \pi_\theta(s). \quad (13)$$

Correspondingly, the “actor” in the actor-critic algorithm can update the parameter θ by

$$\theta \leftarrow \theta + \eta \nabla_a Q_{\pi_\theta}(s, a)|_{a=\pi_\theta(s)} \nabla_\theta \pi_\theta(s). \quad (14)$$

One major issue regarding a deterministic policy is the lack of exploration owing to the determinacy of action selection. To encourage exploration, it is common to perturb the deterministic policy with exploration noises, e.g., adding a Gaussian noise ξ with $a = \pi_\theta(s) + \xi$, for execution.

2) *Modern Actor-Critic Methods*: Although achieving great success in many complex tasks, the actor-critic methods are known to suffer from various problems, such as high variance, slow convergence, local optimum, etc. Therefore, many variants have been developed to improve the performance of actor-critic, and we list some of them as follows.

- *Advantaged Actor-Critic* [2, Chapter 13.4]: The *advantage function*, $A(s, a) = Q_{\pi_\theta}(s, a) - \text{baseline}$, i.e., the Q -function subtracted by a baseline, is introduced to replace $Q_{\pi_\theta}(s, a)$ in the “actor” update, e.g., (9). One common choice for the baseline is an estimate of the state value function $V(s)$. This modification can significantly reduce the variance of the policy gradient estimate without changing the expectation.

- *Asynchronous Actor-Critic* [49] presents an asynchronous variant with parallel training to enhance sample efficiency and training stability. In this method, multiple actors are trained in parallel with different exploration policies, then the global parameters get updated based on all the learning results and synchronized to each actor.

- *Soft Actor-Critic* (SAC) [50] with stochastic policies is an off-policy deep actor-critic algorithm based on the maximum entropy RL framework, which adds an entropy term of the policy $\mathcal{H}(\pi_\theta(\cdot|s_t))$ to objective (1) to encourage exploration.

3) *Trust Region/Proximal Policy Optimization*: To improve the training stability of policy-based RL algorithms, reference [51] proposes the Trust Region Policy Optimization (TRPO) algorithm, which enforces a trust region constraint (15b) that the KL-divergence D_{KL} between the old and new policies should not be greater than a given threshold δ . Specifically, let $\rho(\theta) := \frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)}$ be the probability ratio between the new policy π_θ and the old policy $\pi_{\theta_{old}}$, thus $\rho(\theta_{old}) = 1$. TRPO aims to solve the constrained optimization (15):

$$\max_{\theta} J(\theta) = \mathbb{E}[\rho(\theta)\hat{A}_{\theta_{old}}(s, a)] \quad (15a)$$

$$\text{s.t. } \mathbb{E}[D_{KL}(\pi_\theta || \pi_{\theta_{old}})] \leq \delta \quad (15b)$$

where $\hat{A}_{\theta_{old}}(s, a)$ is an estimation of the advantage function. However, TRPO is relatively complicated to implement. To this end, reference [52] proposes Proximal Policy Optimization (PPO) methods, which achieve the benefits of TRPO with simpler implementation and better empirical sample complexity. Specifically, PPO simplifies the policy optimization as (16):

$$\max_{\theta} \mathbb{E} \left[\min(\rho(\theta)\hat{A}_{\theta_{old}}, \text{clip}(\rho(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_{\theta_{old}}) \right] \quad (16)$$

where ϵ is a hyperparameter, e.g., $\epsilon = 0.2$, and the clip function $\text{clip}(\rho(\theta), 1 - \epsilon, 1 + \epsilon)$ enforces $\rho(\theta)$ to stay within the interval $[1 - \epsilon, 1 + \epsilon]$. The “min” of the clipped and unclipped objective is taken to remove the incentive for moving $\rho(\theta)$ outside of the interval $[1 - \epsilon, 1 + \epsilon]$ and avoid large policy (or θ) update.

4) *Multi-Agent RL*: Many power system control tasks involve the coordination over multiple agents. For example, in frequency regulation, each generator can be treated as an individual agent that makes its own generation decisions, while the frequency dynamics is jointly determined by all power injections. This motivates the *multi-agent RL* framework. Multi-agent RL considers a set \mathcal{N} of agents interacting with the same environment and sharing a common state $s \in \mathcal{S}$. At each time t , each agent $i \in \mathcal{N}$ takes its own action $a_t^i \in \mathcal{A}_i$ given the current state $s_t \in \mathcal{S}$, and receives the reward $r^i(s_t, (a_t^i)_{i \in \mathcal{N}})$, then the system state evolves to s_{t+1} based on $(a_t^i)_{i \in \mathcal{N}}$. Multi-agent RL is an active and challenging research area with many unsolved problems. An overview on related theories and algorithms is provided in [53]. In particular, the decentralized (distributed) multi-agent RL attracts a great deal of attention for power system applications. A popular variant is that each agent i adopts a policy $a^i = \pi_{i, \theta_i}(o^i)$ with its parameter θ_i , which determines the action a^i based on local observation o^i (e.g., local voltage or frequency of bus i). This method allows for decentralized implementation as the policy π_{i, θ_i} for each agent only needs local observations, but it still requires centralized training since the system state transition relies on the actions of all agents. Multi-agent RL methods with distributed training are still under development.

Remark 6. Although the RL techniques above are discussed separately, they can be integrated for a single problem to achieve all the benefits. For instance, one may apply the multi-agent actor-critic framework with deterministic policies, adopt ANNs to parameterize the Q -function and the policy, and use the advantage function for actor update. Accordingly, the resultant algorithm is usually named after the combination of

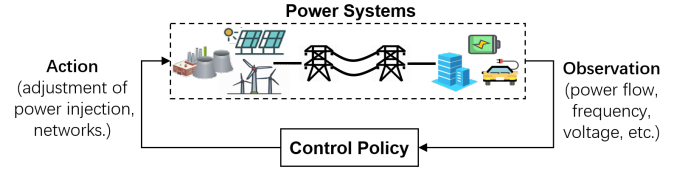


Fig. 5. RL schemes for the control and decision-making in power systems.

key words, e.g., deep deterministic policy gradient (DDPG), asynchronous advantaged actor-critic (A3C), etc.

III. SELECTIVE APPLICATIONS OF RL IN POWER SYSTEMS

Throughout the past decades, tremendous efforts have been devoted to improving the modeling of power systems. Schemes based on (optimal) power flow techniques and precise modeling of various electric facilities are standard for the control and optimization of power systems. However, the large-scale integration of renewable generation and DERs significantly aggravates the complexity, uncertainty, and volatility of power system operation. The network topology, especially for distribution systems, is also changed from time to time due to network reconfiguration, line faults, and other operational factors. It becomes increasingly arduous to procure accurate system models and power injection predictions, challenging the traditional model-based approaches. Hence, model-free RL-based methodology becomes an appealing complement. As illustrated in Figure 5, the RL-based schemes relieve the need for accurate system models and learn control policies based on data collected from real system operation or high-fidelity simulators, whereas the underlying physical models and dynamics are regarded as the unknown environment.

For power systems, frequency level and voltage profile are the two most critical indicators of system operating status, whilst reliable and efficient energy management is the core task. Accordingly, this section focuses on three key applications, i.e., *frequency regulation*, *voltage control*, and *energy management*, and reviews the recent works that apply RL to them. For each of these applications, we summarize the related literature with a table (see Table I, II, III), and explain how to formulate them as RL problems and their solution schemes. The emerging challenges and potential future directions are discussed as well.

Before proceeding, a natural question is why it is necessary to develop new RL-based approaches since traditional tools and existing controllers mostly work “just fine” in real-world power systems. The answer varies from application to application, and we explain some of the main motivations below.

- 1) While traditional methods work well in the current grid, it is envisioned that many traditional schemes may not be sufficient for the future grid with high renewable penetration and user participation. Most existing approaches rely heavily on good knowledge of power system models and have been challenged by many new issues, such as distribution grids that are not well modeled, highly uncertain renewable generation and user behavior, coordination among massive devices, the growing deployment of EVs that are coupled with transportation, and etc.

- 2) The research community has been studying various techniques to tackle these challenges, e.g., distributed control, stochastic optimization, machine learning, etc. Among them, RL is a promising direction to investigate and will play an important role in addressing these challenges due to its data-driven and model-free nature. RL is capable of dealing with highly complex and hard-to-model control problems and can adapt to rapid power fluctuations and topology changes. For example, RL can be used to learn human behavior, especially for demand response and EV charging, and facilitates system operation decisions.
- 3) Lastly, it does not suggest a dichotomy between RL and conventional methods. In fact, RL can be complimentary to existing approaches and improve them in a data-driven way. For instance, policy-based RL algorithms can be integrated to online adjust the parameters of existing controllers for adaptivity and achieve hard-to-model objectives. It is necessary to identify the right application scenarios for RL and use it in an appropriate way. One goal of this paper is to throw light and stimulate such discussions and relevant research.

A. Frequency Regulation

Frequency regulation (FR) is to maintain the power system frequency closely around its nominal value, e.g. 50 Hz or 60 Hz, through balancing power generation and load demand. Conventionally, three generation control mechanisms in a hierarchical structure are implemented at different timescales to achieve fast response and economical efficiency. The primary FR, namely *droop control*, operates locally to eliminate power imbalance at the timescale of few seconds, when the governor adjusts the mechanical power input to the generator around a setpoint and based on the local frequency deviation. The secondary FR, known as *automatic generation control* (AGC), adjusts the setpoints of governors to bring the frequency and tie-line power interchanges back to their nominal values, which is performed in a centralized manner within minutes. The tertiary FR, namely *economic dispatch*, reschedules the unit commitment and restore the secondary control reserves within tens of minutes to hours. See [54] for detailed explanation of the three-level FR architecture. There are a number of recent works [55]–[69] leveraging model-free RL techniques for FR mechanism design, which are summarized in Table I. The main motivations for developing RL-based FR schemes are explained as follows.

- 1) Although the bulk transmission systems have relatively good models on power grids, there may not be accurate models or predictions on the large-scale renewable generation due to the inherent uncertainty and intermittency. As the penetration of renewable generation keeps growing rapidly, new challenges are posed to traditional FR schemes for maintaining nominal frequency in real time.
- 2) FR in the distribution level and demand side, e.g., the PV inverter-based FR, load-side FR, etc., has also attracted a great deal of recent studies. However, distribution systems may not have accurate model information and it is too complex to model massive heterogeneous DERs and load

devices. In such situations, RL methods can be adopted to circumvent the requirement of system model information and learn control policies directly from data.

- 3) With less inertia and fast power fluctuations introduced by the large-scale inverter-based renewable energy resources, power systems become more and more dynamical and volatile. The conventional frequency controllers may not adapt well to the time-varying operational environment [66]. In addition, existing methods have difficulty in coordinating large-scale systems at a fast time scale due to the communication and computation burdens, which limits the overall frequency regulation performance [56]. Hence, (multi-agent) DRL methods are used to develop FR schemes to improve the adaptivity and optimality.

In the follows, we take AGC as the paradigm to illustrate how to apply RL by describing the definition of environment, state and action, the reward design, and the learning of control policies. Then we discuss several key issues in RL-based FR.

1) *Environment, State and Action*: The frequency dynamics in a power network can be expressed as (17):

$$\frac{ds}{dt} = \mathbf{f}(s, \Delta \mathbf{P}_M, \Delta \mathbf{P}_L), \quad (17)$$

where $\mathbf{s} := ((\Delta \omega_i)_{i \in \mathcal{N}}, (\Delta P_{ij})_{ij \in \mathcal{E}})$ denotes the system *state*, including the frequency deviation $\Delta \omega_i$ at each bus i and the power flow deviation ΔP_{ij} from bus i to j (away from the nominal values). $\Delta \mathbf{P}_M := (\Delta P_i^M)_{i \in \mathcal{N}}$, $\Delta \mathbf{P}_L := (\Delta P_i^L)_{i \in \mathcal{N}}$ capture the deviations of generator mechanical power and other power injections, respectively.

The governor-turbine control model [54] of a generator can be formulated as the time differential equation (18):

$$\frac{d\Delta P_i^M}{dt} = g_i(\Delta P_i^M, \Delta \omega_i, P_i^C), \quad i \in \mathcal{N}, \quad (18)$$

where P_i^C is the generation control command. A widely used linearized version of (17) and (18) is provided in Appendix A. However, the real-world frequency dynamics (17) and generation control model (18) are highly nonlinear and complex. This motivates the use of model-free RL methods, since the underlying physical models (17) and (18), together with operational constraints, are simply treated as the *environment* in the RL setting.

When controlling generators for FR, the *action* is defined as the concatenation of the generation control commands with $\mathbf{a} := (P_i^C)_{i \in \mathcal{N}}$. The corresponding action space is continuous in nature but could get discretized in Q -learning based FR schemes [60], [67]. Besides, the continuous-time system dynamics are generally discretized with the discrete-time horizon \mathcal{T} to fit the RL framework, and the time interval Δt depends on the sampling or control period.

We denote $\Delta \mathbf{P}_L$ in (17) as the deviations of other power injections, such as loads (negative power injection), the outputs of renewable energy resources, the charging/discharging power of energy storage systems, and etc. Depending on the actual problem setting, $\Delta \mathbf{P}_L$ could be treated as exogenous states with additional dynamics, or be included in the action \mathbf{a} if these power injections are also controlled for FR [58], [59].

TABLE I
LITERATURE SUMMARY ON LEARNING FOR FREQUENCY REGULATION.

Reference	Problem	State/Action Space	Algorithm	Policy Class	Key Features
Yan et al. 2020 [56]	Multi-area AGC	Continuous	Multi-agent DDPG	ANN	1-Offline centralized learning and decentralized application; 2-Multiply an auto-correlated noise to the actor for exploration; 3-An initialization process is used for ANN training acceleration.
Li et al. 2020 [57]	Single-area AGC	Continuous	Twin delayed DDPG (actor-critic)	ANN	The twin delayed DDPG method is used to improve the exploration process with multiple explorers.
Younesi et al. 2020 [59]	FR in microgrids.	Discrete	Q -learning	ϵ -greedy	Q -learning works as a supervisory controller for PID controllers to improve the online dynamic response.
Khooban et al. 2020 [58]	FR in microgrids.	Continuous	DDPG (actor-critic)	ANN	1-DDPG method works as a supplementary controller for a PID based main controller to improve the online adaptive performance; 2-Add Ornstein-Uhlenbeck process based noises to the actor for exploration.
Chen et al. 2020 [60]	Emergency FR	Discrete	Single/multi-agent Q -learning/DDPG	ϵ -greedy/greedy/ANN	1-Offline learning and online application; 2-The Q -learning/DDPG based controller is used for limited/multiple emergency scenarios.
Abouheaf et al. 2019 [61]	Multi-area AGC	Continuous	Integral RL (actor-critic)	Linear feedback controller	Continuous-time integral-Bellman optimality equation is used.
Wang et al. 2019 [63]	Optimization of activation rules in AGC	Discrete	Multi-objective RL (Q -learning)	Greedy	An constrained optimization model is built to solve for optimal participation factors, where the objective is the combination of multiple Q -functions.
Singh et al. 2017 [67]	Multi-area AGC	Discrete	Q -learning	Stochastic policy	An estimator agent is defined to estimate the frequency bias factor β_i and determine the ACE signal accordingly.

2) *Reward Design*: The design of reward function plays a crucial role in successful RL application. However, there is no general rule to follow, but one principle is to effectively reflect the control goal. For multi-area AGC⁸, it aims to restore the frequency and tie-line power flow to the nominal values after disturbances. Accordingly, the *reward* at time $t \in \mathcal{T}$ can be defined as the minus of frequency deviation and tie-line flow deviation, e.g., the square sum form (19) [56]:

$$r(t) = -\Delta t \cdot \sum_{i \in \mathcal{N}} \left((\beta_i \Delta \omega_i(t))^2 + \left(\sum_{j: ij \in \mathcal{E}} \Delta P_{ij}(t) \right)^2 \right), \quad (19)$$

where β_i is the frequency bias factor. For single-area FR, the goal is to restore the system frequency, thus the term related to tie-line power flow can be removed from (19). Besides, the exponential function [69], absolute value function [57], and other sophisticated reward functions involving the cost of generation change and penalty for large frequency deviation [57], are used as well.

3) *Policy Learning*: Since the system states may not be fully observable in practice, the RL control policy is generally defined as a map $\mathbf{a}(t) = \pi(\mathbf{o}(t))$ from the available measurement observations \mathbf{o} to the action \mathbf{a} instead. The following two steps are critical to effectively learn a good control policy.

- *Select Effective Observation*. Multi-area AGC conventionally operates based on the area control error (ACE) signal, which is defined as $ACE_i = \beta_i \Delta \omega_i + \sum_{j: ij \in \mathcal{E}} \Delta P_{ij}$ with

⁸For multi-area AGC problem, each control area is generally aggregated and characterised by a single governor-turbine model (18). While the control actions for an individual generator within this area are allocated based on its participation factor. Thus each bus i represents an aggregate control area, and ΔP_{ij} is the deviation of tie-line power interchange from area i to area j .

a weighting parameter β_i . Accordingly, the proportional, integral, and derivative (PID) counterparts of the ACE signal, i.e., $(ACE_i(t), \int ACE_i(t) dt, \frac{dACE_i(t)}{dt})$, are adopted as the observation in [56]. Other measurements, such as the power injection deviations $\Delta P_i^M, \Delta P_i^L$, could also be included in the observation [57], [67]. In particular, reference [62] applies the stacked denoising autoencoders to extract compact and useful features from the raw measurement data for FR.

- *Select RL Algorithm*. Both valued-based and policy-based RL algorithms have been applied to FR in power systems. In Q -learning based FR schemes, e.g., [67], the state and action spaces are discretized and the ϵ -greedy policy is used. Recent works [56], [58] employ the DDPG based actor-critic framework to develop the FR schemes, considering continuous action and observation. In addition, multi-agent RL is applied to coordinate multiple control areas or multiple generators in [56], [65], [67], where each agent designs its own control policy $a_i(t) = \pi_i(o_i(t))$ with the local observation o_i . In this way, the resultant algorithms can achieve centralized learning and decentralized implementation.

4) *Discussion*: Based on the existing work described above, we discuss some key observations as follows.

- *Environment Model*. Most references build environment models or simulators to simulate the dynamics and responses of power systems for training and testing their proposed algorithms. These simulators are typically high-fidelity with realistic component models, which are too complex to be useful for direct development and optimization of controllers. Moreover, it is laborious and costly to build and maintain such (dynamical) environment models in practice, and thus they may not be available for many power grids. When such

simulators are unavailable, a potential solution is to train off-policy RL schemes using real system operation data.

- *Safety.* Since FR is vital for power system operation, it necessitates *safe* control policies. Specifically, two requirements need to be met: 1) the closed-loop system dynamics are stable when applying the RL control policies; 2) the physical constraints, such as line thermal limits, are satisfied. However, few existing studies consider the safety issue of applying RL to FR. A recent work [55] proposes to explicitly engineer the ANN structure of DRL to guarantee the frequency stability.

- *Integration with Existing Controllers.* References [58], [59] use the DRL-based controller as a supervisory or supplementary controller to existing PID-based FR controllers, to improve the dynamical adaptivity with baseline performance guarantee. More discussions are provided in Section IV-D.

- *Load-Side Frequency Regulation.* The researches mentioned above focus on controlling generators for FR, while various emerging electric devices, e.g., inverter-based PV units, ubiquitous controllable loads with fast response, are promising complement to generation-frequency control [70], [71]. These are potential FR applications of RL in smart grids.

B. Voltage Control

The goal of voltage control is to keep the voltage magnitudes across the power networks close to the nominal values. Most recent works focus on the voltage control in distribution systems using a variety of control mechanisms [72]–[76]. As the penetration of renewable generation, especially solar panels and wind turbines, deepens in distribution systems, the rapid fluctuations and significant uncertainties of renewable generation pose huge challenges to the voltage control task. Meanwhile, unbalanced power flow, multi-phase device integration, and the lack of accurate network models, further complicate the situation. To this end, a number of studies propose to use model-free RL for voltage control [77]–[89]. We summarize the related work in Table II and present below how to solve the voltage control problem in the RL framework.

1) *Environment, State and Action:* In distribution systems, the controllable devices for voltage control can be classified into slow timescale and fast timescale. Slow timescale devices, such as on-load tap changing transformers (OLTCs), voltage regulators, and capacitor banks, are discretely controlled on an hourly or daily basis. The *states* and control *actions* for them can be defined as

$$\mathbf{s}_{\text{slow}} := ((v_i)_{i \in \mathcal{N}}, (P_{ij}, Q_{ij})_{ij \in \mathcal{E}}, \tau^{\text{TC}}, \tau^{\text{VR}}, \tau^{\text{CB}}), \quad (20a)$$

$$\mathbf{a}_{\text{slow}} := (\Delta \tau^{\text{TC}}, \Delta \tau^{\text{VR}}, \Delta \tau^{\text{CB}}), \quad (20b)$$

where v_i is the voltage magnitude of bus i , and P_{ij}, Q_{ij} are the active and reactive power flows on line ij . $\tau^{\text{TC}}, \tau^{\text{VR}}, \tau^{\text{CB}}$ denote the tap positions of the OLTCs, voltage regulators, and capacitor banks respectively, which are discrete values. $\Delta \tau^{\text{TC}}, \Delta \tau^{\text{VR}}, \Delta \tau^{\text{CB}}$ denote the discrete changes of corresponding tap positions.

The fast timescale devices include inverter-based DERs and static Var compensators (SVCs), whose (active/reactive) power

outputs⁹ can be continuously controlled within seconds. Their *states* and control *actions* can be defined as

$$\mathbf{s}_{\text{fast}} := ((v_i)_{i \in \mathcal{N}}, (P_{ij}, Q_{ij})_{ij \in \mathcal{E}}), \quad (21a)$$

$$\mathbf{a}_{\text{fast}} := (\mathbf{p}^{\text{DER}}, \mathbf{q}^{\text{DER}}, \mathbf{q}^{\text{SVC}}), \quad (21b)$$

where $\mathbf{p}^{\text{DER}}, \mathbf{q}^{\text{DER}}$ collect the continuous active and reactive power outputs of DERs respectively, and \mathbf{q}^{SVC} denotes the reactive power outputs of SVCs.

Since RL methods handle continuous and discrete actions differently, most existing studies only consider either continuous control actions (e.g., \mathbf{q}^{DER}) [85], [88], or discrete control actions (e.g., τ^{CB} and/or τ^{TC}) [82], [83]. While the recent works [84], [91] propose two-timescale/bi-level RL-based voltage control algorithms, taking in account both fast continuous devices and slow discrete devices. In the follows, we uniformly use \mathbf{s} and \mathbf{a} to denote the state and action.

Given the definitions of state and action, the system dynamics that depict the *environment* can be formulated as

$$\mathbf{s}(t+1) = \mathbf{f}(\mathbf{s}(t), \mathbf{a}(t), \mathbf{p}^{\text{ex}}(t), \mathbf{q}^{\text{ex}}(t)), \quad (22)$$

where $\mathbf{p}^{\text{ex}}, \mathbf{q}^{\text{ex}}$ denote the exogenous active power and reactive power injections to the grid, including load demands and other generations. The transition function \mathbf{f} captures the tap position evolution and the power flow equations, which can be very complex or even unknown in reality. The exogenous injections $\mathbf{p}^{\text{ex}}, \mathbf{q}^{\text{ex}}$ include the uncontrollable renewable generations that are difficult to predict as well. Similar to FR, the dynamics model (22) is not required in RL algorithms.

2) *Reward Design:* The goal of voltage control is to design control policies such that the voltage magnitudes are close to the nominal value, which we denote as 1 per unit (p.u.). Accordingly, the reward function is typically in the form of penalization on the voltage deviation from 1 p.u. For example, the reward can simply be (23), e.g., in [83], [84], [88],

$$r(\mathbf{s}, \mathbf{a}) = - \sum_{i \in \mathcal{N}} (v_i - 1)^2, \quad (23)$$

where the negative sign indicates the smaller the voltage violation, the higher the reward. An alternative way is to set the reward to be negative (e.g. -1) when the voltage is outside a tolerance range (e.g. $\pm 5\%$ of the nominal value), and positive (e.g. $+1$) when inside the range, e.g. in [87]. Furthermore, the reward can incorporate the operation cost of the controllable devices (e.g. switching cost of discrete devices) and the power loss [86], as well as some sophisticated definitions [85].

3) *RL Algorithms:* Both value-based and policy-based RL algorithms have been applied for voltage control:

- *Value-Based RL.* A number of works [82]–[84], [87] adopt value-based algorithms, such as DQN and LSPI, to learn the optimal Q -function with function approximation, typically using ANNs [84], [87] or radial basis functions [83]. Based on the learned Q -function, the control policy is chosen as the greedy policy, i.e., $\mathbf{a}(t) = \arg \max_{\tilde{\mathbf{a}} \in \mathcal{A}} Q(\mathbf{s}(t), \tilde{\mathbf{a}})$. Two

⁹Due to the comparable magnitudes of line resistance and reactance in distribution networks, the conditions for active-reactive power decoupling are no longer met. Thus *active power* outputs also play a role in voltage control, and alternating current (AC) power flow models are generally needed.

TABLE II
LITERATURE SUMMARY ON LEARNING FOR VOLTAGE CONTROL.

Reference	Control Scheme	State/Action Space	Algorithm	Policy Class	Key Features
Gao et. al. 2021 [77]	Voltage regulator, capacitor banks, OLTC	Hybrid/discrete	Consensus multi-agent DRL	ANN	1- The maximum entropy method is used to encourage exploration; 2- A consensus multi-agent RL algorithm is developed, which enables distributed control and efficient communication.
Sun et. al. 2021 [78]	PV reactive power control	Hybrid/continuous	Multi-agent DDPG	ANN	A voltage sensitivity based DDPG method is proposed, which analytically computes the gradient of value function to action rather than using the critic ANN.
Zhang et. al. 2021 [79]	Smart inverters, voltage regulators, and capacitors	Continuous/discrete	Multi-agent DQN	ϵ -greedy	1- Both the network loss and voltage violation are considered in the reward definition; 2- Multi-agent DQN is used to enhance the scalability of the algorithm.
Mukherjee et. al. 2021 [80]	Load shedding	Continuous	Hierarchical DRL	LSTM	1- A hierarchical multi-agent RL algorithm with two levels is developed to accelerate learning; 2- Augmented random search is used to solve for optimal policies.
Kou et. al. 2020 [81]	Reactive power control	Continuous	DDPG	ANN	A safety layer is formed on top of the actor network to ensure safe exploration, which predicts the state change and prevents the violation of constraints.
Xu et. al. 2020 [83]	Set OLTC tap positions	Hybrid/discrete	LSPI (batch RL)	Greedy	1- “Virtual” sample generation is used for better exploration; 2- Adopt a multi-agent trick to handle scalability; 3-Use Radial basis function as features for state.
Yang et. al. 2020 [84]	On-off switch of capacitors	Hybrid/discrete	DQN	Greedy	A two-time scale scheme with power injections determined via traditional OPF in fast time scale, and the switching of capacitors determined via DQN in slow time scale.
Wang et. al. 2020 [85]	Generator voltage setpoints	Continuous	Multi-agent DDPG	ANN	Adopt a competitive (game) formulation with specially designed reward for each agent.
Wang et. al. 2020 [86]	Set tap/on-off positions	Hybrid/Discrete	Constrained SAC	ANN	1- Model the voltage violation as constraint using the Constrained MDP framework; 2- Reward is defined as the negative of power loss and switching cost.
Duan et. al. 2020 [87]	Generator voltage setpoints	Hybrid	DQN/DDPG	Decaying ϵ -greedy/ANN	1-DQN is used for discrete action and DDPG is used for continuous action; 2- Decaying ϵ -greedy policy is employed in DQN to encourage exploration.
Cao et. al. 2020 [88]	PV reactive power control	Continuous	Multi-agent DDPG	ANN	The attention neural network is used to develop the critic to enhance the algorithm scalability.
Liu et. al. 2020 [89]	Reactive power control	Continuous	Adversarial RL [90]/SAC	Stochastic policy	1-A two-stage RL method is proposed to improve the online safety and efficiency via offline pre-training; 2- Adversarial SAC is used to make the online application robust to the transfer gap.

limitations of the greedy policy include 1) the action selection depends on the state of the entire system, which hinders distributed implementation; 2) it is usually not suitable for continuous action space since the maximization is not easy to compute, especially when complex function approximation, e.g., with ANNs, is adopted.

• *Policy-Based RL.* Compared with value-based RL methods, the voltage control schemes based on actor-critic algorithms, e.g., [85]–[88], are more flexible, which can accommodate both continuous and discrete actions and enable distributed implementation. For example, a parameterized deterministic policy class $q_i^{\text{DER}} = \pi_{i,\theta_i}(\mathbf{o}_i^{\text{DER}})$ is employed for each DER device i , which determines the reactive power output q_i^{DER} based on local observation $\mathbf{o}_i^{\text{DER}}$ with parameter θ_i . The policy class π_{i,θ_i} is often parameterized using ANNs. Then some actor-critic methods, e.g., multi-agent DDPG, are used to optimize parameter θ_i , where a centralized critic learns the Q -function with ANN approximation and each individual DER device performs the policy gradient update as the actor.

4) *Discussion:* As the network scale and the number of controllable devices increase, the size of the state/action space grow exponentially, which poses serious challenges in learning the Q -function. It is worth mentioning that reference [83] proposes a special trick that defines different Q -functions for different actions, which addresses the scalability issue under its special problem formulation. Besides, to learn the Q -function for a given policy, on-policy RL methods, such as actor-critic, need to implement the policy and collect sample data. This could be problematic since the policy is not necessarily safe, and thus the implementation on real-world power systems may be catastrophic. One remedy is to train the policy on high-fidelity simulators. Reference [83] proposes a novel method to generate “virtual” sample data for a certain policy, based on the data collected from implementing another safe policy. More discussions on safety are provided in Section IV-A.

C. Energy Management

Energy management is an advanced application that utilizes information flow to manage power flow and maintain power balance in a reliable and efficient manner. To this end, energy management systems (EMSs) are developed for electric power control centers to monitor, control, and optimize the system operation. With the assistance of the supervisory control and data acquisition (SCADA) system, the EMS for transmission systems is technically mature. However, for many sub-regional power systems, such as medium/low-voltage distribution grids and microgrids, EMS is still under development due to the integration of various DER facilities and lack of metering units. Moreover, a EMS family [92] with a hierarchical structure is necessitated to facilitate different levels of energy management, including grid-level EMS, EMS for coordinating a cluster of DERs, home EMS (HEMS), etc.

In practice, there exist significant uncertainties that result from unknown models and parameters of power networks and DER facilities, uncertain user behaviors and weather conditions, etc. Hence, many recent researches adopt (D)RL techniques to develop data-driven EMS. A summary of the related literature is provided in Table III. In the follows, we first introduce the models of DERs and adjustable loads, then review the RL based schemes for different levels of energy management problems.

1) *State, Action, and Environment*: We present the action, state and environment models for several typical DER facilities, buildings, and residential loads.

• *Distributed Energy Resources*: For compact expression, we consider a bundle of several typical DERs, including a dispatchable photovoltaics (PV) unit, a battery, an electric vehicle (EV), and a diesel generator (DG). The *action* at time $t \in \mathcal{T}$ is defined as

$$\mathbf{a}^{\text{DER}}(t) := (p^{\text{PV}}(t), p^{\text{Bat}}(t), p^{\text{EV}}(t), p^{\text{DG}}(t)). \quad (24)$$

Here, $p^{\text{PV}}, p^{\text{Bat}}, p^{\text{EV}}, p^{\text{DG}}$ are the power outputs of PV, battery, EV, and DG respectively, which are continuous in nature. $p^{\text{Bat}}, p^{\text{EV}}$ can be either positive (discharging) or negative (charging). The DER *state* at time t can be defined as

$$\mathbf{s}^{\text{DER}}(t) := (\bar{p}^{\text{PV}}(t), E^{\text{Bat}}(t), E^{\text{EV}}(t), \mathbf{x}^{\text{EV}}(t)), \quad (25)$$

where \bar{p}^{PV} is the maximal PV generation power determined by the solar irradiance. The PV output p^{PV} can be adjusted within the interval $[0, \bar{p}^{\text{PV}}]$, and $p^{\text{PV}} = \bar{p}^{\text{PV}}$ when the PV unit operates in the maximum power point tracking (MPPT) mode. $E^{\text{Bat}}, E^{\text{EV}}$ denote the associated state of charge (SOC) levels. \mathbf{x}^{EV} captures other related states of EV, e.g. current location (at home or outside), travel plan, etc.

• *Building HVAC*: Buildings account for a large share of the total energy usage, about half of which is consumed by the heating, ventilation, and air conditioning (HVAC) systems [93]. Smartly scheduling HVAC operation has huge potential to save energy cost, but the building climate dynamics is intrinsically hard to model and affected by various environmental factors. Generally, a building is divided into multiple thermal zones, and the *action* at time t is defined as

$$\mathbf{a}^{\text{HVAC}}(t) := (T_c(t), T_s(t), (m^i(t))_{i \in \mathcal{N}}), \quad (26)$$

where T_c and T_s are the conditioned air temperature and the supply air temperature, respectively. m^i is the supply air flow rate at zone $i \in \mathcal{N}$. The choice of states is subtle, since many exogenous factors may affect the indoor climate. A typical definition of the *HVAC state* is

$$\mathbf{s}^{\text{HVAC}}(t) := (T_{\text{out}}(t), (T_{\text{in}}^i(t), h^i(t), e^i(t))_{i \in \mathcal{N}}), \quad (27)$$

where T_{out} and T_{in}^i are the outside temperature and indoor temperature of zone i ; h^i and e^i are the humidity and occupancy rate of zone i , respectively. Besides, the solar irradiance, the carbon dioxide concentration and other environmental factors may also be included in the state \mathbf{s}^{HVAC} .

• *Residential Loads*: Residential demand response (DR) [94] that motivates changes in electric consumption by end users in response to time-varying electricity price or incentive payments attracts considerable recent attention. The domestic electric appliances are classified as 1) *non-adjustable loads*, e.g., computers, refrigerators, which are critical and must be satisfied; 2) *adjustable loads*, e.g., air conditioner, washing machine, whose operating power or time can be tuned. The *action* for an adjustable load i at time $t \in \mathcal{T}$ is defined as

$$\mathbf{a}_i^{\text{L}}(t) := (z_i^{\text{L}}(t), p_i^{\text{L}}(t)), \quad i \in \mathcal{N}_L, \quad (28)$$

where binary $z_i^{\text{L}} \in \{0, 1\}$ denotes whether switching the on/off mode (equal 1) or keeping unchanged (equal 0). p_i^{L} is the power consumption of load i , which can be adjusted either discretely or continuously depending on the load characteristics. The operational *state* of load i can be defined as

$$\mathbf{s}_i^{\text{L}}(t) := (\alpha_i^{\text{L}}(t), \mathbf{x}_i^{\text{L}}(t)), \quad i \in \mathcal{N}_L, \quad (29)$$

where binary α_i^{L} equals 0 for the off status and 1 for the on status. \mathbf{x}_i^{L} collects other related states of load i . For example, the indoor and outdoor temperatures are contained in \mathbf{x}_i^{L} if load i is an air condition; \mathbf{x}_i^{L} captures the task progress and remaining time to the deadline for a washing machine load.

• *Other System States*: In addition to the operational states above, there are some critical *system states* for EMS, e.g.,

$$\mathbf{s}^{\text{Sys}}(t) := (t, \ell(t - K_p : t + K_f), \mathbf{v}(t), \mathbf{P}(t), \dots), \quad (30)$$

including the current time t , electricity price ℓ (from past K_p time steps to future K_f time predictions) [95], voltage profile $\mathbf{v} := (v_i)_{i \in \mathcal{N}}$, power flow $\mathbf{P} := (P_{ij})_{ij \in \mathcal{E}}$, and etc.

The state definitions (25) (27) (29) only contain the present status at time t , while the past values and future predictions are often included in the state as well to capture the temporal patterns. In addition, the previous actions may also be considered as components of the state, e.g., adding $\mathbf{a}^{\text{DER}}(t-1)$ to the state $\mathbf{s}^{\text{DER}}(t)$. For different energy management problems, the state \mathbf{s} and action \mathbf{a} are determined accordingly by selecting and combining the definitions in (but not limited to) (24)-(30). And the *environment* model is given by

$$\mathbf{s}(t+1) = \mathbf{f}(\mathbf{s}(t), \mathbf{a}(t), \mathbf{u}^{\text{ex}}(t)), \quad (31)$$

where \mathbf{u}^{ex} captures other related exogenous factors. We note that the RL models presented above are used as illustrative examples, while one needs to formulate its own models to fit the specific applications.

2) *Energy Management Applications*: Energy management indeed covers a broad range of sub-topics, including integrated energy systems (IESs), grid-level power dispatch, management of DERs, building HVAC control, and HEMS, etc. We present these sub-topics in a hierarchical order as follows, and summarize the basic attributes and key features of representative references in Table III.

- *Integrated Energy Systems* [96], also referred as multi-energy systems, couple the power grids with heat networks and gas networks, to accommodate renewable energy and enhance the overall energy efficiency and flexibility. Reference [97] proposes a DDPG-based real-time control strategy to manage a residential multi-energy system, where DERs, heat pumps, gas boilers, and thermal energy storage are controlled to minimize the total operational cost. In [98], the management of IESs with integrated demand response is modeled as a Stackelberg game, and an actor-critic scheme is developed for the energy provider to adjust pricing and power dispatching strategies to cope with unknown privacy parameters of users. Extensive case studies are conducted in [99] to compare the performance of a twin delayed DDPG scheme against a benchmark linear model-predictive-control method, which empirically show that RL is a viable optimal control technique for IES management and could outperform conventional tools.

- *Grid-Level Power Dispatch* aims to schedule the power outputs of generators and DERs to optimize the operating cost of the entire grid, while satisfying the operational constraints of electric facilities and the network. Optimal power flow (OPF) is the fundamental tool of traditional power dispatch schemes. Several recent works [100]–[104] propose DRL-based methods to solve the OPF problem in order to achieve fast solution and address the absence of accurate grid models. Most existing references [105]–[110] focus on the power dispatch in distribution grids or microgrids. In [108], a model-based DRL algorithm is developed to online schedule a residential microgrid, and Monte-Carlo tree search is adopted to make the optimal decisions. Reference [110] proposes a cooperative RL algorithm for distributed economic dispatch in microgrids, where a diffusion strategy is used to coordinate the actions of DERs.

- *Device-Level Energy Management* focuses on the optimal control of DER devices and adjustable loads, such as EV, energy storage system, HVAC, and residential electric appliances, which generally aims to minimize the total energy cost under time-varying electricity price. In [111]–[113], various RL techniques are studied to design EV charging policies to deal with the randomness in the arrival and departure time of an EV. See [114] for a review on RL-based EV charging management systems. References [115]–[117] adopt DQN and DDPG to learn the charging/discharging strategy for controlling battery systems considering unknown degradation models. In terms of building HVAC control, there are multiple uncertainty factors such as random zone occupancy, unknown thermal dynamics models, uncertain outdoor temperature and electricity price, etc., while the thermal comfort and air quality comfort need to be guaranteed. Therefore, a number of studies [118]–[122] leverage DRL for HVAC system control. In [95], [123]–[125], DRL-based HEMS is developed to optimally

schedule household electric appliances, taking into account the resident's preference as well as uncertain electricity price and weather conditions.

3) *Discussion*: Some key issues are discussed as follows.

- *Physical Constraints*. For practical energy management, there are many physical constraints, e.g., the state of charge limits for batteries and EVs, that should be satisfied when taking control actions. Reference [110] formulates the constraint violation as a penalty term in the reward function, which is in the form of a logarithmic barrier function. Reference [112] builds a constrained MDP problem to take the physical constraints into account and solves the problem with the constrained policy optimization method [113]. These methods actually impose the constraints in a “soft” manner, where the implemented actions may still have the chance to violate the constraints. More discussions are provided in Section IV-A.

- *Hybrid of Discrete and Continuous State/Action*. Energy management often involves the control of a hybrid of discrete devices and continuous devices, while the basic RL methods only focus on handling either discrete or continuous actions. Some *Q*-learning based work [123] simply discretizes the continuous action space to fit the algorithm framework. Reference [124] proposes a ANN based stochastic policy to handle both discrete and continuous actions, which is a combination of the Bernoulli policy for on/off switch action and the Gaussian policy for continuous action.

D. Other Applications

In addition to the three critical applications above, other applications of RL in power systems include electricity market [129], [130], network reconfiguration [131], service restoration [132], [133], emergency control [134], maximum power point tracking [135], [136], cyber security [137], [138], maintenance scheduling [139], protective relay control [140], electric vehicle charging navigation [141], demand response customer selection [142], power flexibility aggregation [143], and etc.

IV. CHALLENGES AND PERSPECTIVES

In this section, we present three critical challenges of using RL in power system applications, i.e., safety, scalability, and data. Several potential future directions are then discussed.

A. Safety

Power systems are vital infrastructures of modern societies, thus it must ensure that any controllers applied are safe, in the sense that they do not drive the power system operational states to violate crucial physical constraints, or cause instability or reliability issues. Regarding RL based control schemes, there are the following two aspects of safety concern:

- 1) *Guarantee that the learning process is safe* (also referred as *safe exploration*). For this issue, off-policy RL methods [83] are more desired, where the training data are generated from existing controllers that are known to be safe. In contrast, it remains an open question for on-policy RL to guarantee safe exploration. Some attempts [144]–[147] propose safe on-policy exploration schemes based on Lyapunov criterion and

TABLE III
LITERATURE SUMMARY ON LEARNING FOR ENERGY MANAGEMENT.

Reference	Problem	State/Action Space	Algorithm	Policy Class	Key Features
Ye et al. 2020 [97]	IES management	Hybrid	DDPG (actor-critic)	Gaussian policy	The prioritized experience replay method is used to enhance the sampling efficiency of the experience replay mechanism.
Xu et al. 2021 [126]	IES management	Discrete	Q -learning	Stochastic policy	A RL-based differential evolution algorithm is developed to solve the complex IES scheduling issue.
Yan et al. 2020 [100]	Optimal power flow	Continuous	DDPG	ANN	Rather than using the critic network, the deterministic gradient of a Lagrangian-based value function is derived analytically.
Zhou et al. 2020 [103]	Optimal power flow	Continuous	Proximal policy optimization	Gaussian policy	Imitation learning is adopted to generate initial weights for ANNs and proximal policy optimization is used to train a DRL agent for fast OPF solution.
Hao et al. 2021 [107]	Microgrid power dispatch	Continuous	Hierarchical RL	Two knowledge rule-based policies	1- Hierarchical RL is used to reduce complexity and improve learning efficiency; 2- Incorporated with domain knowledge, it avoids baseline violation and additional learning beyond feasible action space.
Lin et al. 2020 [105]	Power dispatch	Continuous	Soft A3C	Gaussian policy	The edge computing technique is employed to accelerate the computation and communication in a cloud environment.
Zhang et al. 2020 [106]	Distribution power dispatch	Continuous	Fitted Q -iteration	ϵ -greedy	The Q -function is parameterized by polynomial approximation and optimized using a regularized recursive least square method with a forgetting factor.
Wan et al. 2019 [111]	EV charging scheduling	Continuous/discrete	Double DQN	ϵ -greedy	A representation network is constructed to extract features from the electricity price.
Li et al. 2020 [112]	EV charging scheduling	Continuous	Constrained policy optimization [113]	Gaussian policy	A constrained MDP is formulated to schedule the charging of a EV, considering the charging constraints.
Silva et al. 2020 [127]	EV charging scheduling	Discrete	Multi-agent Q -learning	ϵ -greedy	Use multi-agent multi-objective RL to mode the EV charging coordination with the W -learning method.
Bui et al. 2020 [115]	Battery management	Hybrid/discrete	Double DQN	ϵ -greedy	To mitigate the overestimation problem, double DQN with a primary network for action selection and a target network is used.
Cao et al. 2020 [116]	Battery management	Hybrid/discrete	Double DQN	Greedy	A hybrid CNN and LSTM model is adopted to predict the future electricity price.
Yu et al. 2021 [120]	Building HVAC	Continuous	Multi-actor-attention-critic	Stochastic Policy	A scalable HVAC control algorithm is proposed to solve the Markov game based on multi-agent DRL with attention mechanism.
Gao et al. 2020 [119]	Building HVAC	Continuous	DDPG	ANN	A feed-forward ANN with Bayesian regularization is built for predicting thermal comfort.
Mocanu et al. 2019 [121]	Building energy management	Hybrid	DPG and DQN	ANN	Both DPG and DQN are implemented for building energy control and their performances are compared numerically.
Xu et al. 2020 [123]	HEMS	Continuous/discrete	Multi-agent Q -learning	ϵ -greedy	Use extreme learning machine based ANNs to predict future PV output and electricity price.
Li et al. 2020 [124]	HEMS	Hybrid	Trust region policy optimization	ANNs based stochastic policy	A policy network determines both discrete actions (on/off switch with Bernoulli policy) and continuous actions (power control with Gaussian policy).
Chen et al. 2021 [128]	Residential load control	Continuous	Thompson sampling	Stochastic policy	Logistic regression is employed to predict customers' opt-out behaviors in demand response and Thompson sampling is used for online learning.

Gaussian process. The basic idea is to construct a certain safety region, and special actions are taken to drive the state back once approaching the boundary of this safety region. See [148] for a comprehensive survey on safe RL. However, almost all the existing works train their RL control policies only based on high-fidelity power system simulators, and it is plausible that the safe exploration problem is circumvented. Nevertheless, one can argue that there might be a substantial gap between the simulator and the real-world system, leading to the failure of generalization in real implementation. To this end, a possible

remedy is to employ the robust (adversarial) RL methods [89], [90], [149] in simulator-based policy training.

2) *Guarantee that the final learned control policy is safe.* It is generally hard to verify whether a policy is safe or its generated actions can respect physical operational constraints. Some common methods to deal with constraints include 1) formulating the constraint violation as a penalty term to the reward; 2) training the control policy based on constrained MDP [86], [112]. Specifically, the second way aims to learn an optimal policy π^* that maximizes the expected total return

$J(\pi)$ and is subject to a budget constraint:

$$\pi^* \in \arg \max_{\pi} J(\pi), \quad \text{s.t. } J^c(\pi) \leq d, \quad (32)$$

where $J^c(\pi)$ is the expected total cost and d is the budget. By defining the physical constraint violation as certain costs in $J^c(\pi)$, (32) imposes safety requirements to some degrees. Typical approaches to solve the constrained MDP problem (32) include the Lagrangian methods [86], [150], constrained policy update rules [113], and etc.

Besides, constrained RL is an active research area that deals with safety and constraint issues. Two types of constraints, i.e., soft constraints and hard constraints, are generally considered in literature. The typical ways to handle soft constraints include 1) using barrier functions or penalty functions to integrate the constraints to the reward function [151]; 2) modeling in the form of a chance constraint (i.e., $\mathbb{P}(\text{constraint violation}) \leq \epsilon$, where ϵ is a probability threshold) [152], [153] or a long-term constraint (such as the constraint in model (32)) [154], [155]. In terms of hard constraints, the predominant approach is to take conservative actions to ensure that the hard constraints on state and action are satisfied at all time, despite the uncertainty in the problem [156]. However, such schemes usually lead to significant conservativeness and may not work well when the system becomes complex, as in the case of power systems.

B. Scalability

In most of the existing work, it is observed that only small-scale power systems with a few decision-making agents are tested using simulations. To the best of our knowledge, no real-life implementation of RL control schemes has been reported yet. A crucial limitation for RL in large-scale multi-agent systems, such as power systems, is the scalability issue, since the state and action spaces expand dramatically as the number of agents increases, which is known as the “curse of dimensionality”. The multi-agent RL and function approximation techniques are useful to improve the scalability, but they are still under development with many limitations. For instance, there is limited provable guarantee on how well Q -function can be approximated with ANNs, and it is unclear whether it works for real-size power grids. Moreover, even though an individual policy based on local observations can be used for each agent, most existing multi-agent RL methods still need centralized learning among all the agents, as the Q -function depends on the global state and the actions of all the agents. To this end, scalable actor-critic methods that enable distributed learning for networked multi-agent systems are proposed in [157], [158], which leverage the local dependency properties to find (near-)optimal localized policies. Besides, some application-specific approximation methods can be utilized to design scalable RL algorithms. For example, reference [83] develops a scalable LSPI-based voltage control scheme, which uses the trick that sequentially learns a separate approximate Q -functions for each component of the action, whilst the other components are assumed to behave greedily according to their own approximate Q -function.

C. Data

1) *Data Quantity and Quality*: Evaluating the amount of data that are needed for training a good policy, namely *sample complexity*, is an important issue and active research area in the RL community. For classical RL algorithms, such as Q -learning, the sample complexity depends on the size of the state and action spaces of the problem; generally, the larger the state and action spaces are, the more data are needed to find a near optimal policy [25], [159]. For modern RL methods that are commonly used in power systems, such as DQN and actor-critic, the sample complexity also depends on the complexity of the function class adopted to approximate the Q -function, as well as the intrinsic approximation error of the function class [46], [160]. In addition, data quality is one of the critical factors affecting the learning efficiency. Real measurement and operational data of power grids suffer from various issues, such as missing data, outlier data, noise data, etc., thus a pre-processing on raw data is needed. Theoretically, larger variance in noisy observations typically leads to higher sample complexity for achieving a certain level of accuracy. Despite the successful application in simulations, theoretic understanding on the sample complexity of modern RL algorithms is limited and deserves more studies. Moreover, many power system applications employ the multi-agent training methods with partial observation, which further complicates the theoretical analysis.

2) *Data Availability*: Almost all references reviewed above assume that high-fidelity simulators or accurate environment models are available to simulate the system dynamics and response, which are the sources of sample data for training and testing RL policies. However, when such simulators are unavailable, data availability becomes an issue for the application of on-policy RL algorithms. One potential solution is to construct training samples from the consecutive fine-grained system operational data and adopt off-policy RL methods to learn control policies. However, off-policy RL usually is less efficient and has slower convergence than on-policy RL, which is an ongoing research in the RL community.

3) *Standardized Dataset and Testbed*: It is observed that different synthetic test systems and test cases are used in the power literature to simulate and test the proposed RL-based algorithms, and many implementation details and codes are not provided. Hence, it is necessary to develop benchmark datasets and authoritative testbeds for power system applications, in order to standardize the testing of RL algorithms and facilitate fair performance comparison.

4) *Big Data Techniques*: The big data in smart grids can benefit the application of data-driven RL in multiple ways [161], which include 1) measurement data from SCADA, PMUs, AMI, and other advanced metering devices, 2) electricity market pricing and bidding data, 3) equipment monitoring, control, maintenance, and management data, 4) meteorological data, etc. Specifically, big data mining techniques for knowledge discovery can be adopted to detect special events, determine effective observations, and identify critical latent states. Pattern extraction from massive datasets can be utilized to classify and cluster similar events, agents and user behaviors, to improve the data efficiency and scalability of RL algorithms.

D. Future Directions

Regarding the challenges in applying RL to power systems, we present several potential future directions as below.

1) *Integrate model-free and model-based methods*: The practical power system operation is not a black box and does have useful model information to work with. Purely model-free approaches may be too radical to exploit available information and suffer from their own limitations, such as the safety and scalability issues discussed above. Since existing model-based methods have already been well studied in theory and applied in industry with acceptable performance, one promising future direction is to *combine model-based and model-free methods* for complementarity and achieve the advantages of both. For example, model-based methods could work as a warm-start, the nominal model, or be used to identify critical features for model-free methods. Model-free methods could be used to coordinate and adjust the parameters of incumbent model-based controllers to improve their adaptivity. In [162], the potential integration of model-based and model-free methods is summarized as three types: implementing them in serial, in parallel, or embedding one as an inner module in the other. Although there is limited work on this subject so far, the integration of model-free RL techniques with existing model-based control schemes is envisioned to be an important future research.

2) *Exploit Suitable RL Variants*: RL is a fundamental and vibrant research field that is attracting a great deal of attention. New advances in RL algorithms appear frequently. In addition to DRL, multi-agent RL, and robust RL mentioned above, a wide range of branches in the RL field, such as transfer RL [163], meta RL [164], federated RL [165], inverse RL [37], integral RL [14], Bayesian RL [166], hierarchical RL [167], interpretable RL [168], etc., can be utilized to improve the learning efficiency and address specific problems in appropriate application scenarios. For instance, transfer RL can be employed to transplant the well-trained policies for one task to another similar task, so that it does not have to learn from scratch and thus can enhance the training efficiency.

3) *Leverage Domain Knowledge and Problem Structures*: The naive application of existing RL algorithms may encounter many troubles in practice. In addition to algorithmic advances, *leveraging domain knowledge and exploiting application-specific structures to design tailored RL algorithms* are necessary to achieve superior performance. Specifically, domain knowledge and empirical evidences can be used to guide the definition of state and reward, the initialization of the policy, as well as the selection of basic RL algorithms. For example, as presented above, the area control error (ACE) signal is adopted as the state when applying RL for frequency regulation. Besides, the specific problem structures are useful in determining the policy class, approximation function class, hyperparameters, etc., to improve training efficiency and provide performance guarantee. For example, reference [55] leverages two special properties of the frequency control problem to design the policy network with a particular structure, so that the resultant RL controllers have stability guarantee.

4) *Satisfy Practical Requirements*: To achieve practical implementation in power systems, the following concrete

requirements on RL-based methods need to be met:

- As discussed above, the safety and scalability issues of RL based methods need to be addressed.
- RL based algorithms should be *robust* to the noises and failures in measurement, communication, computation, and actuation, to ensure reliable operation.
- To be used with confidence, RL based methods need to be *interpretable* and have theoretical performance guarantee.
- Since RL generally requires a large amount data from multi-stakeholders, the data privacy should be preserved.
- As power systems generally operate under normal conditions, it remains an unsolved problem to make sure that the RL control policies learned from real system data have sufficient exploration and can perform well in extreme scenarios.
- Since RL based approaches heavily rely on information flow, the cyber security should be ensured under various malicious cyber attacks.
- Existing RL based algorithms mostly take tens of thousands of iterations to converge, which suggests that the *training efficiency* needs to be improved.
- Necessary computing resources, communications infrastructure and technology need to be deployed or upgraded to support the application of RL schemes.

V. CONCLUSION

Although a number of works have been devoted to applying RL to the power system field, many key problems remain unsolved and there is still a substantial distance from practical implementation. On the one hand, this subject is new and still under development and needs much more studies. On the other hand, it is time to step back and rethink the advantages and limitations of applying RL to power systems (the world's most complex and vital engineered systems) and figure out where and when to apply RL. In fact, RL is not envisioned to completely replace existing model-based methods but a viable alternative in specific tasks. For instance, when the models are too complex to be useful, or when the problems are intrinsically hard to model, such as the human-in-loop control (e.g., in demand response), RL and other data-driven methods are promising. Identifying the right application scenarios for RL and using it in an appropriate way are highly expected.

APPENDIX A SYSTEM FREQUENCY DYNAMICS

According to [54], [67], [71], the system frequency dynamics (17) can be linearized as (33), with the generator swing dynamics (33a) and the power flow dynamics (33b).

$$\Delta \dot{\omega}_i = -\frac{1}{M_i}(D_i \Delta \omega_i - \Delta P_i^M + \Delta P_i^L + \sum_{j:ij \in \mathcal{E}} \Delta P_{ij}), i \in \mathcal{N} \quad (33a)$$

$$\Delta \dot{P}_{ij} = B_{ij}(\Delta \omega_i - \Delta \omega_j), ij \in \mathcal{E} \quad (33b)$$

where M_i, D_i, B_{ij} denote the generator inertia, damping coefficient, and synchronization coefficient, respectively. Besides, the governor-turbine control model (18) for a generator can be

simplified as (34), including the turbine dynamics (34a) and the governor dynamics (34b):

$$\Delta \dot{P}_i^M = -\frac{1}{T_i^{\text{tur}}}(\Delta P_i^M - \Delta P_i^G) \quad (34a)$$

$$\Delta \dot{P}_i^G = -\frac{1}{T_i^{\text{gov}}}(\frac{1}{R_i}\Delta\omega_i + \Delta P_i^G - P_i^C) \quad (34b)$$

where ΔP_i^G is turbine valve position deviation, and P_i^C is the generation control command. $T_i^{\text{tur}}, T_i^{\text{gov}}$ denote the turbine and governor time constants, and R_i is the droop coefficient.

REFERENCES

- [1] V. C. Gungor, D. Sahin, T. Kocak, S. Ergut, C. Buccella, C. Cecati, and G. P. Hancke, "Smart grid technologies: Communication technologies and standards," *IEEE Trans. Ind. Informat.*, vol. 7, no. 4, pp. 529–539, 2011.
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2018.
- [3] Y. Li, "Deep reinforcement learning: An overview," *arXiv preprint arXiv:1701.07274*, 2017.
- [4] D. Silver, A. Huang, C. J. Maddison *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [5] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *The Int. J. of Robot. Res.*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [6] A. E. Sallab, M. Abdou, E. Perot, and S. Yogamani, "Deep reinforcement learning framework for autonomous driving," *Electronic Imaging*, vol. 2017, no. 19, pp. 70–76, 2017.
- [7] Y. Zhao, D. Zeng, M. A. Socinski, and M. R. Kosorok, "Reinforcement learning strategies for clinical trials in nonsmall cell lung cancer," *Biometrics*, vol. 67, no. 4, pp. 1422–1433, 2011.
- [8] Z. Zhang, D. Zhang, and R. C. Qiu, "Deep reinforcement learning for power system applications: An overview," *CSEE J. Power Energy Syst.*, vol. 6, no. 1, pp. 213–225, 2019.
- [9] M. Glavic, "(Deep) reinforcement learning for electric power system control and related problems: A short review and perspectives," *Annual Reviews in Control*, vol. 48, pp. 22–35, 2019.
- [10] D. Cao, W. Hu, J. Zhao, G. Zhang, B. Zhang, Z. Liu, Z. Chen, and F. Blaabjerg, "Reinforcement learning and its applications in modern power and energy systems: A review," *J. Mod. Power Syst. Clean Energy*, vol. 8, no. 6, pp. 1029–1042, 2020.
- [11] T. Yang, L. Zhao, W. Li, and A. Y. Zomaya, "Reinforcement learning in sustainable energy and electric systems: A survey," *Annual Reviews in Control*, vol. 49, pp. 145–163, 2020.
- [12] T. M. Moerland, J. Broekens, and C. M. Jonker, "Model-based reinforcement learning: A survey," *arXiv preprint arXiv:2006.16712*, 2020.
- [13] D. Bertsekas, *Dynamic Programming and Optimal Control*. Athena Scientific Belmont, MA, 2012.
- [14] H. Modares, F. L. Lewis, and M.-B. Naghibi-Sistani, "Integral reinforcement learning and experience replay for adaptive optimal control of partially-unknown constrained-input continuous-time systems," *Automatica*, vol. 50, no. 1, pp. 193–202, 2014.
- [15] E. Lecarpentier and E. Rachelson, "Non-stationary markov decision processes, a worst-case approach using model-based reinforcement learning, extended version," *arXiv preprint arXiv:1904.10090*, 2019.
- [16] W. C. Cheung, D. Simchi-Levi, and R. Zhu, "Reinforcement learning for non-stationary markov decision processes: The blessing of (more) optimism," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2020, pp. 1843–1854.
- [17] T. Jaksch, R. Ortner, and P. Auer, "Near-optimal regret bounds for reinforcement learning," *J. Mach. Learn. Res.*, vol. 11, no. 4, 2010.
- [18] D. Russo, B. Van Roy, A. Kazerouni, I. Osband, and Z. Wen, "A tutorial on thompson sampling," *arXiv preprint arXiv:1707.02038*, 2017.
- [19] J. N. Tsitsiklis, "Asynchronous stochastic approximation and Q-learning," *Machine Learning*, vol. 16, no. 3, pp. 185–202, 1994.
- [20] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan, "Optimality and approximation with policy gradient methods in markov decision processes," in *Conf. Learning Theory*. PMLR, 2020, pp. 64–66.
- [21] J. N. Tsitsiklis and B. Van Roy, "An analysis of temporal-difference learning with function approximation," *IEEE Trans. Autom. Control*, vol. 42, no. 5, pp. 674–690, 1997.
- [22] H. Robbins and S. Monro, "A stochastic approximation method," *Ann. Math. Stat.*, pp. 400–407, 1951.
- [23] R. Srikant and L. Ying, "Finite-time error bounds for linear stochastic approximation and TD learning," *arXiv preprint arXiv:1902.00923*, 2019.
- [24] C. J. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3–4, pp. 279–292, 1992.
- [25] G. Qu and A. Wierman, "Finite-time analysis of asynchronous stochastic approximation and q-learning," in *Conf. Learning Theory*. PMLR, 2020, pp. 3185–3205.
- [26] G. A. Rummery and M. Niranjan, *On-line Q-Learning Using Connectionist Systems*. University of Cambridge, Department of Engineering Cambridge, UK, 1994.
- [27] T. Degris, M. White, and R. S. Sutton, "Off-policy actor-critic," *arXiv preprint arXiv:1205.4839*, 2012.
- [28] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Adv. Neural Inf. Process. Syst.*, 2000, pp. 1057–1063.
- [29] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *Proc. Int. Conf. Machine Learning*, 2014, pp. 387–395.
- [30] S. Lange, T. Gabel, and M. Riedmiller, "Batch reinforcement learning," in *Reinforcement learning*. Springer, 2012, pp. 45–73.
- [31] D. Ernst, P. Geurts, and L. Wehenkel, "Tree-based batch mode reinforcement learning," *J. Mach. Learn. Res.*, vol. 6, pp. 503–556, Apr. 2005.
- [32] M. G. Lagoudakis and R. Parr, "Least-squares policy iteration," *J. Mach. Learn. Res.*, vol. 4, pp. 1107–1149, Dec. 2003.
- [33] R. Agarwal, D. Schuurmans, and M. Norouzi, "An optimistic perspective on offline reinforcement learning," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2020, pp. 104–114.
- [34] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," *arXiv preprint arXiv:2005.01643*, 2020.
- [35] D. Krueger, J. Leike, O. Evans, and J. Salvatier, "Active reinforcement learning: Observing rewards at a cost," *arXiv preprint arXiv:2011.06709*, 2020.
- [36] C. Daniel, M. Viering, J. Metz, O. Kroemer, and J. Peters, "Active reward learning," in *Robotics: Science and systems*, vol. 98, 2014.
- [37] A. Y. Ng, S. J. Russell *et al.*, "Algorithms for inverse reinforcement learning," in *ICML*, vol. 1, 2000, p. 2.
- [38] J. Fu, K. Luo, and S. Levine, "Learning robust rewards with adversarial inverse reinforcement learning," *arXiv preprint arXiv:1710.11248*, 2017.
- [39] H. Wang and B. Raj, "On the origin of deep learning," *arXiv preprint arXiv:1702.07800*, 2017.
- [40] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [41] R. Sun, "Optimization for deep learning: theory and algorithms," *arXiv preprint arXiv:1912.08957*, 2019.
- [42] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [43] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [44] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [45] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [46] J. Fan, Z. Wang, Y. Xie, and Z. Yang, "A theoretical analysis of deep q-learning," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 486–489.
- [47] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proc. AAAI Conf. on Artif. Intel.*, vol. 30, no. 1, 2016.
- [48] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2016, pp. 1995–2003.
- [49] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 1928–1937.
- [50] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," *arXiv preprint arXiv:1801.01290*, 2018.

- [51] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 1889–1897.
- [52] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [53] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *arXiv preprint arXiv:1911.10635*, 2019.
- [54] H. Bevrani and T. Hiyama, *Intelligent Automatic Generation Control*. CRC press, 2017.
- [55] W. Cui and B. Zhang, "Reinforcement learning for optimal frequency control: A lyapunov approach," *arXiv preprint arXiv:2009.05654*, 2020.
- [56] Z. Yan and Y. Xu, "A multi-agent deep reinforcement learning method for cooperative load frequency control of a multi-area power system," *IEEE Trans. Power Syst.*, vol. 35, no. 6, Nov. 2020.
- [57] J. Li and T. Yu, "Deep reinforcement learning based multi-objective integrated automatic generation control for multiple continuous power disturbances," *IEEE Access*, vol. 8, pp. 156 839–156 850, 2020.
- [58] M. H. Khooban and M. Gheisarnajad, "A novel deep reinforcement learning controller based type-II fuzzy system: Frequency regulation in microgrids," *IEEE Trans. Emerg. Topics Comput. Intell.*, pp. 1–11, 2020.
- [59] A. Younesi, H. Shayeghi, and P. Siano, "Assessing the use of reinforcement learning for integrated voltage/frequency control in AC microgrids," *Energies*, vol. 13, no. 5, 2020.
- [60] C. Chen, M. Cui, F. F. Li, S. Yin, and X. Wang, "Model-free emergency frequency control based on reinforcement learning," *IEEE Trans. Ind. Informat.*, pp. 1–1, 2020.
- [61] M. Abouheaf, Q. Gueaieb, and A. Sharaf, "Load frequency regulation for multi-area power system using integral reinforcement learning," *IET Gener. Transm. Distrib.*, vol. 13, no. 19, pp. 4311–4323, 2019.
- [62] Z. Yan and Y. Xu, "Data-driven load frequency control for stochastic power systems: A deep reinforcement learning method with continuous action search," *IEEE Trans. Power Syst.*, vol. 34, no. 2, pp. 1653–1656, Mar. 2019.
- [63] H. Wang, Z. Lei, X. Zhang, J. Peng, and H. Jiang, "Multiobjective reinforcement learning-based intelligent approach for optimization of activation rules in automatic generation control," *IEEE Access*, vol. 7, pp. 17 480–17 492, Feb. 2019.
- [64] M. Adibi and J. van der Woude, "A reinforcement learning approach for frequency control of inverted-based microgrids," *IFAC-PapersOnLine*, vol. 52, no. 4, pp. 111–116, 2019.
- [65] L. Xi, J. Chen, Y. Huang, Y. Xu, L. Liu, Y. Zhou, and Y. Li, "Smart generation control based on multi-agent reinforcement learning with the idea of the time tunnel," *Energy*, vol. 153, pp. 977–987, Jun. 2018.
- [66] L. Yin, T. Yu, L. Zhou, L. Huang, X. Zhang, and B. Zheng, "Artificial emotional reinforcement learning for automatic generation control of large-scale interconnected power grids," *IET Gener. Transm. Distrib.*, vol. 11, no. 9, pp. 2305–2313, 2017.
- [67] V. P. Singh, N. Kishor, and P. Samuel, "Distributed multi-agent system-based load frequency control for multi-area power system in smart grid," *IEEE Trans. Ind. Electron.*, vol. 64, no. 6, pp. 5151–5160, Jun. 2017.
- [68] T. Yu, H. Z. Wang, B. Zhou, K. W. Chan, and J. Tang, "Multi-agent correlated equilibrium Q(λ) learning for coordinated smart generation control of interconnected power grids," *IEEE Trans. Power Syst.*, vol. 30, no. 4, pp. 1669–1679, Jul. 2015.
- [69] S. Rozada, D. Apostolopoulou, and E. Alonso, "Load frequency control: A deep multi-agent reinforcement learning approach," in *Proc. IEEE Power Energy Soc. Gen. Meeting*, Montreal, Canada, Aug. 2020.
- [70] E. Mallada, C. Zhao, and S. Low, "Optimal load-side control for frequency regulation in smart grids," *IEEE Trans. Autom. Control*, vol. 62, no. 12, pp. 6294–6309, 2017.
- [71] X. Chen, C. Zhao, and N. Li, "Distributed automatic load frequency control with optimality in power systems," *IEEE Control Netw. Syst.*, 2020.
- [72] H. Sun, Q. Guo, and et al., "Review of challenges and research opportunities for voltage control in smart grids," *IEEE Trans. Power Syst.*, vol. 34, no. 4, pp. 2790–2801, 2019.
- [73] K. E. Antoniadou-Plytaria, I. N. Kouveliotis-Lysikatos, P. S. Georgilakis, and N. D. Hatziaargyriou, "Distributed and decentralized voltage control of smart distribution networks: Models, methods, and future research," *IEEE Trans. Smart Grid*, vol. 8, no. 6, pp. 2999–3008, 2017.
- [74] G. Qu and N. Li, "Optimal distributed feedback voltage control under limited reactive power," *IEEE Trans. Power Syst.*, vol. 35, no. 1, pp. 315–331, 2019.
- [75] S. Magnússon, G. Qu, and N. Li, "Distributed optimal voltage control with asynchronous and delayed communication," *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3469–3482, 2020.
- [76] S. M. Mohiuddin and J. Qi, "Droop-free distributed control for ac microgrids with precisely regulated voltage variance and admissible voltage profile guarantees," *IEEE Trans. Smart Grid*, vol. 11, no. 3, pp. 1956–1967, 2020.
- [77] Y. Gao, W. Wang, and N. Yu, "Consensus multi-agent reinforcement learning for volt-var control in power distribution networks," *IEEE Trans. Smart Grid*, pp. 1–1, 2021.
- [78] X. Sun and J. Qiu, "Two-stage volt/var control in active distribution networks with multi-agent deep reinforcement learning method," *IEEE Trans. Smart Grid*, pp. 1–1, 2021.
- [79] Y. Zhang, X. Wang, J. Wang, and Y. Zhang, "Deep reinforcement learning based volt-var optimization in smart distribution systems," *IEEE Trans. Smart Grid*, vol. 12, no. 1, pp. 361–371, 2021.
- [80] S. Mukherjee, R. Huang, Q. Huang, T. L. Vu, and T. Yin, "Scalable voltage control using structure-driven hierarchical deep reinforcement learning," *arXiv preprint arXiv:2102.00077*, 2021.
- [81] P. Kou, D. Liang, C. Wang, Z. Wu, and L. Gao, "Safe deep reinforcement learning-based constrained optimal control scheme for active distribution networks," *Applied Energy*, vol. 264, p. 114772, 2020.
- [82] J. Vlachogiannis and N. Hatziaargyriou, "Reinforcement learning for reactive power control," *IEEE Trans. Power Syst.*, vol. 19, no. 3, pp. 1317–1325, Aug. 2004.
- [83] H. Xu, A. D. Domínguez-García, and P. W. Sauer, "Optimal tap setting of voltage regulation transformers using batch reinforcement learning," *IEEE Trans. Power Syst.*, vol. 35, no. 3, pp. 1990–2001, May 2020.
- [84] Q. Yang, G. Wang, A. Sadeghi, G. B. Giannakis, and J. Sun, "Two-timescale voltage control in distribution grids using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 3, pp. 2313–2323, May 2020.
- [85] S. Wang, J. Duan, D. Shi, C. Xu, H. Li, R. Diao, and Z. Wang, "A data-driven multi-agent autonomous voltage control framework using deep reinforcement learning," *IEEE Trans. Power Syst.*, 2020.
- [86] W. Wang, N. Yu, Y. Gao, and J. Shi, "Safe off-policy deep reinforcement learning algorithm for Volt-VAR control in power distribution systems," *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3008–3018, Jul. 2020.
- [87] J. Duan, D. Shi, R. Diao, H. Li, Z. Wang, B. Zhang, D. Bian, and Z. Yi, "Deep-reinforcement-learning-based autonomous voltage control for power grid operations," *IEEE Trans. Power Syst.*, vol. 35, no. 1, pp. 814–817, Jan. 2020.
- [88] D. Cao, W. Hu, J. Zhao, Q. Huang, Z. Chen, and F. Blaabjerg, "A multi-agent deep reinforcement learning based voltage regulation using coordinated PV inverters," *IEEE Trans. Power Syst.*, 2020.
- [89] H. Liu and W. Wu, "Two-stage deep reinforcement learning for inverter-based volt-var control in active distribution networks," *IEEE Trans. Smart Grid*, 2020.
- [90] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," *arXiv preprint arXiv:1703.02702*, 2017.
- [91] H. Liu and W. Wu, "Bi-level off-policy reinforcement learning for volt/var control involving continuous and discrete devices," *arXiv preprint arXiv:2104.05902*, 2021.
- [92] H. Sun, B. Zhang, W. Wu, and Q. Guo, "Family of energy management system for smart grid," in *Proc. 3rd IEEE PES Innov. Smart Grid Technol. Int. Conf. Exhibit.*, Berlin, Germany, 2012, pp. 1–5.
- [93] Office of Energy Efficiency & Renewable Energy (EERE). 2011 buildings energy data book. [Online]. Available: <https://catalog.data.gov/dataset/buildings-energy-data-book>.
- [94] U.S. Department of Energy, "Benefit of demand response in electricity market and recommendations for achieving them," Feb. 2006.
- [95] R. Lu, S. H. Hong, and M. Yu, "Demand response for home energy management using reinforcement learning and artificial neural network," *IEEE Trans. Smart Grid*, vol. 10, no. 6, pp. 6629–6639, 2019.
- [96] W. Huang, N. Zhang, Y. Cheng, J. Yang, Y. Wang, and C. Kang, "Multienergy networks analytics: standardized modeling, optimization, and low carbon analysis," *Proceedings of the IEEE*, vol. 108, no. 9, pp. 1411–1436, 2020.
- [97] Y. Ye, D. Qiu, X. Wu, G. Strbac, and J. Ward, "Model-free real-time autonomous control for a residential multi-energy system using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3068–3082, 2020.
- [98] Y. Wang, Z. Yang, L. Dong, S. Huang, and W. Zhou, "Energy management of integrated energy system based on stackelberg game

- and deep reinforcement learning," in *Proc. IEEE 4th Conf. Energy Inter. and Energy System Integ.*, 2020, pp. 2645–2651.
- [99] G. Ceusters, R. C. Rodríguez, A. B. García, R. Franke, G. Deconinck, L. Helsen, A. Nowé, M. Messagie, and L. R. Camargo, "Model-predictive control and reinforcement learning in multi-energy system case studies," *arXiv preprint arXiv:2104.09785*, 2021.
- [100] Z. Yan and Y. Xu, "Real-time optimal power flow: A lagrangian based deep reinforcement learning approach," *IEEE Trans. Power Syst.*, vol. 35, no. 4, pp. 3270–3273, 2020.
- [101] C. Jiang, Z. Li, J. H. Zheng, Q. H. Wu, and X. Shang, "Two-level area-load modelling for opf of power system using reinforcement learning," *IET Gener., Transm. Distrib.*, vol. 13, no. 18, pp. 4141–4149, 2019.
- [102] J. H. Woo, L. Wu, J.-B. Park, and J. H. Roh, "Real-time optimal power flow using twin delayed deep deterministic policy gradient algorithm," *IEEE Access*, vol. 8, pp. 213 611–213 618, 2020.
- [103] Y. Zhou, B. Zhang, C. Xu, T. Lan, R. Diao, D. Shi, Z. Wang, and W.-J. Lee, "A data-driven method for fast ac optimal power flow solutions via deep reinforcement learning," *J. Mod. Power Syst. Clean Energy*, vol. 8, no. 6, pp. 1128–1139, 2020.
- [104] D. Cao, W. Hu, X. Xu, Q. Wu, Q. Huang, Z. Chen, and F. Blaabjerg, "Deep reinforcement learning based approach for optimal power flow of distribution networks embedded with renewable energy and storage devices," *J. Mod. Power Syst. Clean Energy*, pp. 1–10, 2021.
- [105] L. Lin, X. Guan, Y. Peng, N. Wang, S. Maharjan, and T. Ohtsuki, "Deep reinforcement learning for economic dispatch of virtual power plant in internet of energy," *IEEE Internet of Things J.*, vol. 7, no. 7, pp. 6288–6301, 2020.
- [106] Q. Zhang, K. Dehghanpour, Z. Wang, and Q. Huang, "A learning-based power management method for networked microgrids under incomplete information," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1193–1204, 2020.
- [107] R. Hao, T. Lu, Q. Ai, and H. He, "Distributed online dispatch for microgrids using hierarchical reinforcement learning embedded with operation knowledge," *IEEE Trans. Power Syst.*, pp. 1–1, 2021.
- [108] H. Shuai and H. He, "Online scheduling of a residential microgrid via monte-carlo tree search and a learned model," *IEEE Trans. Smart Grid*, vol. 12, no. 2, pp. 1073–1087, 2021.
- [109] Y. Du and F. Li, "Intelligent multi-microgrid energy management based on deep neural network and model-free reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1066–1076, 2020.
- [110] W. Liu, P. Zhuang, H. Liang, J. Peng, and Z. Huang, "Distributed economic dispatch in microgrids based on cooperative reinforcement learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 6, pp. 2192–2203, 2018.
- [111] Z. Wan, H. Li, H. He, and D. Prokhorov, "Model-free real-time EV charging scheduling based on deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 10, no. 5, pp. 5246–5257, Sep. 2019.
- [112] H. Li, Z. Wan, and H. He, "Constrained EV charging scheduling based on safe deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 3, pp. 2427–2439, May 2020.
- [113] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *Proc. of the 34th Int. Conf. Mach. Learn.*, 2017, p. 22–31.
- [114] H. M. Abdullah, A. Gastli, and L. Ben-Brahim, "Reinforcement learning based ev charging management systems—a review," *IEEE Access*, vol. 9, pp. 41 506–41 531, 2021.
- [115] V.-H. Bui, A. Hussain, and H.-M. Kim, "Double deep Q -learning-based distributed operation of battery energy storage system considering uncertainties," *IEEE Trans. Smart Grid*, vol. 11, no. 1, pp. 457–469, Jan. 2020.
- [116] J. Cao, D. Harrold, Z. Fan, T. Morstyn, D. Healey, and K. Li, "Deep reinforcement learning-based energy storage arbitrage with accurate lithium-ion battery degradation model," *IEEE Trans. Smart Grid*, vol. 11, no. 5, pp. 4513–4521, 2020.
- [117] F. Sanchez Gorostiza and F. M. Gonzalez-Longatt, "Deep reinforcement learning-based controller for soc management of multi-electrical energy storage system," *IEEE Trans. Smart Grid*, vol. 11, no. 6, pp. 5039–5050, 2020.
- [118] T. Wei, Y. Wang, and Q. Zhu, "Deep reinforcement learning for building HVAC control," in *Proc. of the 54th Annual Design Autom. Conf.*, 2017, pp. 1–6.
- [119] G. Gao, J. Li, and Y. Wen, "Deepcomfort: Energy-efficient thermal comfort control in buildings via reinforcement learning," *IEEE Internet of Things J.*, vol. 7, no. 9, pp. 8472–8484, 2020.
- [120] L. Yu, Y. Sun, Z. Xu, C. Shen, D. Yue, T. Jiang, and X. Guan, "Multi-agent deep reinforcement learning for HVAC control in commercial buildings," *IEEE Trans. Smart Grid*, vol. 12, no. 1, pp. 407–419, 2021.
- [121] E. Mocanu, D. C. Mocanu, P. H. Nguyen, A. Liotta, M. E. Webber, M. Gibescu, and J. G. Sloatweg, "On-line building energy optimization using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 10, no. 4, pp. 3698–3708, 2019.
- [122] X. Zhang, D. Biagioni, M. Cai, P. Graf, and S. Rahman, "An edge-cloud integrated solution for buildings demand response using reinforcement learning," *IEEE Trans. Smart Grid*, vol. 12, no. 1, pp. 420–431, 2021.
- [123] X. Xu, Y. Jia, Y. Xu, Z. Xu, S. Chai, and C. S. Lai, "A multi-agent reinforcement learning based data-driven method for home energy management," *IEEE Trans. Smart Grid*, 2020.
- [124] H. Li, Z. Wan, and H. He, "Real-time residential demand response," *IEEE Trans. Smart Grid*, vol. 11, no. 5, pp. 4144–4154, Sep. 2020.
- [125] F. Alfaverh, M. Denai, and Y. Sun, "Demand response strategy based on reinforcement learning and fuzzy reasoning for home energy management," *IEEE Access*, vol. 8, pp. 39 310–39 321, 2020.
- [126] Z. Xu, G. Han, L. Liu, M. Martínez-García, and Z. Wang, "Multi-energy scheduling of an industrial integrated energy system by reinforcement learning based differential evolution," *IEEE Trans. Green Commun. Netw.*, pp. 1–1, 2021.
- [127] F. L. Silva, C. E. H. Nishida, D. M. Roijers, and A. H. R. Costa, "Coordination of electric vehicle charging through multiagent reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 3, pp. 2347–2356, 2020.
- [128] X. Chen, Y. Li, J. Shimada, and N. Li, "Online learning and distributed control for residential demand response," *IEEE Trans. Smart Grid*, pp. 1–1, 2021.
- [129] Y. Ye, D. Qiu, J. Li, and G. Strbac, "Multi-period and multi-spatial equilibrium analysis in imperfect electricity markets: A novel multi-agent deep reinforcement learning approach," *IEEE Access*, vol. 7, pp. 130 515–130 529, 2019.
- [130] Y. Ye, D. Qiu, M. Sun, D. Papadaskalopoulos, and G. Strbac, "Deep reinforcement learning for strategic bidding in electricity markets," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1343–1355, 2020.
- [131] Y. Gao, W. Wang, J. Shi, and N. Yu, "Batch-constrained reinforcement learning for dynamic distribution network reconfiguration," *IEEE Trans. Smart Grid*, 2020.
- [132] L. R. Ferreira, A. R. Aoki, and G. Lambert-Torres, "A reinforcement learning approach to solve service restoration and load management simultaneously for distribution networks," *IEEE Access*, vol. 7, pp. 145 978–145 987, 2019.
- [133] D. Ye, M. Zhang, and D. Sutanto, "A hybrid multiagent framework with q-learning for power grid systems restoration," *IEEE Trans. Power Syst.*, vol. 26, no. 4, pp. 2434–2441, 2011.
- [134] Q. Huang, R. Huang, W. Hao, J. Tan, R. Fan, and Z. Huang, "Adaptive power system emergency control using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1171–1182, 2020.
- [135] C. Wei, Z. Zhang, W. Qiao, and L. Qu, "Reinforcement-learning-based intelligent maximum power point tracking control for wind energy conversion systems," *IEEE Trans. Ind. Electron.*, vol. 62, no. 10, pp. 6360–6370, 2015.
- [136] A. Kushwaha, M. Gopal, and B. Singh, "Q-learning based maximum power extraction for wind energy conversion system with variable wind speed," *IEEE Trans. Energy Convers.*, vol. 35, no. 3, pp. 1160–1170, 2020.
- [137] D. An, Q. Yang, W. Liu, and Y. Zhang, "Defending against data integrity attacks in smart grid: A deep reinforcement learning-based approach," *IEEE Access*, vol. 7, pp. 110 835–110 845, 2019.
- [138] Y. Chen, S. Huang, F. Liu, Z. Wang, and X. Sun, "Evaluation of reinforcement learning-based false data injection attack to automatic voltage control," *IEEE Trans. Smart Grid*, vol. 10, no. 2, pp. 2158–2169, 2019.
- [139] Y. Shang, W. Wu, J. Liao, J. Guo, J. Su, W. Liu, and Y. Huang, "Stochastic maintenance schedules of active distribution networks based on monte-carlo tree search," *IEEE Trans. Power Syst.*, vol. 35, no. 5, pp. 3940–3952, 2020.
- [140] D. Wu, X. Zheng, D. Kalathil, and L. Xie, "Nested reinforcement learning based control for protective relays in power distribution systems," in *Proc. IEEE 58th Conf. Decision and Control*, 2019, pp. 1925–1930.
- [141] T. Qian, C. Shao, X. Wang, and M. Shahidehpour, "Deep reinforcement learning for EV charging navigation by coordinating smart grid and intelligent transportation system," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1714–1723, 2020.
- [142] X. Chen, Y. Nie, and N. Li, "Online residential demand response via contextual multi-armed bandits," *IEEE Contr. Syst. Lett.*, vol. 5, no. 2, pp. 433–438, 2021.

- [143] T. Li, B. Sun, Y. Chen, Z. Ye, S. H. Low, and A. Wierman, "Real-time aggregate flexibility via reinforcement learning," *arXiv preprint arXiv:2012.11261*, 2020.
- [144] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees," in *Adv. Neural Inf. Process. Syst.*, 2017, pp. 908–918.
- [145] T. Koller, F. Berkenkamp, M. Turchetta, and A. Krause, "Learning-based model predictive control for safe exploration," in *Proc. IEEE Conf. on Decision and Control*, Miami Beach, FL, 2018, pp. 6059–6066.
- [146] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning," in *Adv. Neural Inf. Process. Syst.*, 2018, pp. 8092–8101.
- [147] J. Fan and W. Li, "Safety-guided deep reinforcement learning via online gaussian process estimation," *arXiv preprint arXiv:1903.02526*, 2019.
- [148] J. Garcia and F. Fernández, "A comprehensive survey on safe reinforcement learning," *J. Mach. Learn. Res.*, vol. 16, no. 1, pp. 1437–1480, 2015.
- [149] J. Morimoto and K. Doya, "Robust reinforcement learning," *Neural Comput.*, vol. 17, no. 2, pp. 335–359, 2005.
- [150] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," *arXiv preprint arXiv:1805.11074*, 2018.
- [151] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick, "End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks," in *Proc. AAAI Conf. Artif. Intel.*, vol. 33, no. 01, 2019, pp. 3387–3395.
- [152] Z. Qin, Y. Chen, and C. Fan, "Density constrained reinforcement learning," in *Inter. Conf. Mach. Learn.* PMLR, 2021, pp. 8682–8692.
- [153] K. P. Wabersich and M. N. Zeilinger, "Linear model predictive safety certification for learning-based control," in *2018 IEEE Conf. Decision and Control (CDC)*. IEEE, 2018, pp. 7130–7135.
- [154] E. Altman, *Constrained Markov decision processes*. CRC Press, 1999, vol. 7.
- [155] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," *arXiv preprint arXiv:1705.10528*, 2017.
- [156] S. Dean, S. Tu, N. Matni, and B. Recht, "Safely learning to control the constrained linear quadratic regulator," in *2019 Amer. Contr. Conf. (ACC)*. IEEE, 2019, pp. 5582–5588.
- [157] G. Qu, A. Wierman, and N. Li, "Scalable reinforcement learning of localized policies for multi-agent networked systems," *Proc. Mach. Learn. Res.*, vol. 1, p. 38, 2020.
- [158] Y. Lin, G. Qu, L. Huang, and A. Wierman, "Multi-agent reinforcement learning in time-varying networked systems," *arXiv preprint arXiv:2006.06555*, 2020.
- [159] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen, "Sample complexity of asynchronous q-learning: Sharper analysis and variance reduction," *arXiv preprint arXiv:2006.03041*, 2020.
- [160] H. Kumar, A. Koppel, and A. Ribeiro, "On the sample complexity of actor-critic method for reinforcement learning with function approximation," *arXiv preprint arXiv:1910.08412*, 2019.
- [161] M. Ghorbanian, S. H. Dolatabadi, and P. Siano, "Big data issues in smart grids: A survey," *IEEE Systems Journal*, vol. 13, no. 4, pp. 4158–4168, 2019.
- [162] Q. Wang, F. Li, Y. Tang, and Y. Xu, "Integrating model-driven and data-driven methods for power system frequency stability assessment and control," *IEEE Trans. Power Syst.*, vol. 34, no. 6, pp. 4557–4568, Nov. 2019.
- [163] E. Parisotto, J. L. Ba, and R. Salakhutdinov, "Actor-mimic: Deep multitask and transfer reinforcement learning," *arXiv preprint arXiv:1511.06342*, 2015.
- [164] J. X. Wang, Z. Kurth-Nelson, D. Tirumala, and et al., "Learning to reinforcement learn," *arXiv preprint arXiv:1611.05763*, 2016.
- [165] H. H. Zhuo, W. Feng, Q. Xu, Q. Yang, and Y. Lin, "Federated reinforcement learning," *arXiv preprint arXiv:1901.08277*, vol. 1, 2019.
- [166] M. Ghavamzadeh, S. Mannor, J. Pineau, and A. Tamar, "Bayesian reinforcement learning: A survey," *arXiv preprint arXiv:1609.04436*, 2016.
- [167] A. G. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning," *Discrete Event Dynamic Syst.*, vol. 13, no. 1, pp. 41–77, 2003.
- [168] A. Verma, V. Murali, R. Singh, P. Kohli, and S. Chaudhuri, "Programmatically interpretable reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 5045–5054.

Xin Chen received the double B.S. degrees in engineering physics and economics and the master's degree in electrical engineering from Tsinghua University, Beijing, China, in 2015 and 2017, respectively. He is currently pursuing the Ph.D. degree in electrical engineering with Harvard University, USA. He was a recipient of the Best Conference Paper Award in the IEEE PES General Meeting in 2016 and the Best Student Paper Award Finalist in the IEEE Conference on Control Technology and Applications in 2018. His research interests lie in reinforcement learning, distributed optimization and control of networked systems, with emphasis on power and energy systems.

Guannan Qu received his B.S. degree in Electrical Engineering from Tsinghua University in Beijing, China in 2014, and his Ph.D. in Applied Mathematics from Harvard University in Cambridge, MA in 2019. Since 2019 he has been a CMI and Resnick postdoctoral scholar in the Department of Computing and Mathematical Sciences at California Institute of Technology. He is the recipient of Simoudis Discovery Award, PIMCO Fellowship, and AI4Science Fellowship at Caltech, and Best Student Paper Reward from IEEE SmartGridComm. His research interest lies in control, optimization, and machine/reinforcement learning with applications to power systems, multi-agent systems, smart city, etc.

Yujie Tang received his B.S. degree in Electronics Engineering from Tsinghua University in 2013 and his Ph.D. degree in Electrical Engineering from the California Institute of Technology in 2019. He is a postdoctoral fellow in the School of Engineering and Applied Sciences at Harvard University. His research interest lies in distributed and real-time/online optimization and their applications in cyber-physical networks.

Steven Low (F' 2008) is the F. J. Gilloon Professor of the Department of Computing & Mathematical Sciences and the Department of Electrical Engineering at Caltech. Before that, he was with AT&T Bell Laboratories, Murray Hill, NJ, and the University of Melbourne, Australia. He has held honorary/chaired professorship Australia, China and Taiwan. He was a co-recipient of IEEE best paper awards and is a Fellow of both IEEE (2008) and ACM (2020). He was a member of the Networking and Information Technology Technical Advisory Group for the US President's Council of Advisors on Science and Technology (PCAST) in 2006. He received his B.S. from Cornell and PhD from Berkeley, both in EE. His research areas include smart grid, cyber-physical systems, network architecture, and energy-efficient networking.

Na Li is a Gordon McKay professor in Electrical Engineering and Applied Mathematics at Harvard University. She received her Bachelor degree in Mathematics from Zhejiang University in 2007 and Ph.D. degree in Control and Dynamical systems from California Institute of Technology in 2013. She was a postdoctoral associate at Massachusetts Institute of Technology 2013–2014. Her research lies in control, learning, and optimization of networked systems, including theory development, algorithm design, and applications to real-world cyber-physical societal system in particular power systems. She received NSF career award (2016), AFSOR Young Investigator Award (2017), ONR Young Investigator Award (2019), Donald P. Eckman Award (2019), McDonald Mentoring Award (2020), along with some other awards.